

Motion learning in robotics

Karel Zimmermann

<http://cmp.felk.cvut.cz/~zimmerk/>



Vision for Robotics and Autonomous Systems

<https://cyber.felk.cvut.cz/vras/>



Center for Machine Perception

<https://cmp.felk.cvut.cz>



Department for Cybernetics
Faculty of Electrical Engineering
Czech Technical University in Prague



Tasks often formalised as MDP

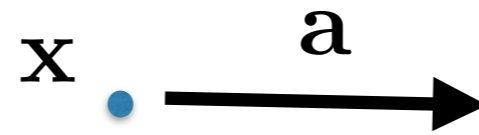
States: $\mathbf{x} \in \mathcal{R}^n$

\mathbf{x} ●



Both tasks formalised as reinforcement learning problems

States: $\mathbf{x} \in \mathcal{R}^n$



Actions: $\mathbf{a} \in \mathcal{R}^m$

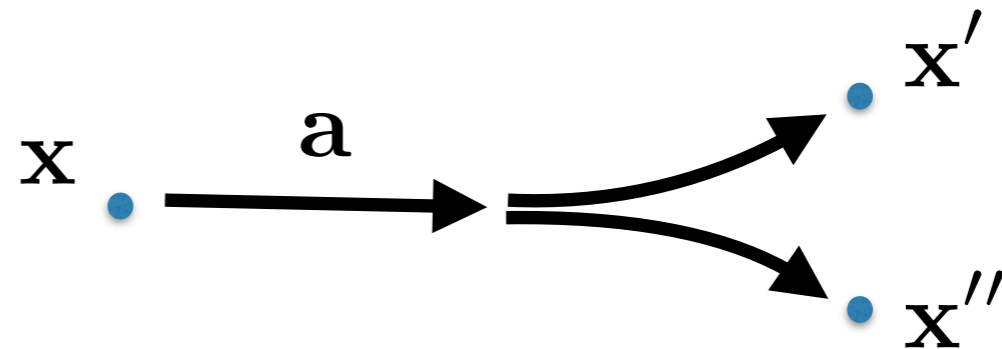


Both tasks formalised as reinforcement learning problems

States: $\mathbf{x} \in \mathcal{R}^n$

Actions: $\mathbf{a} \in \mathcal{R}^m$

Model: $p(\mathbf{x}' | \mathbf{x}, \mathbf{a})$



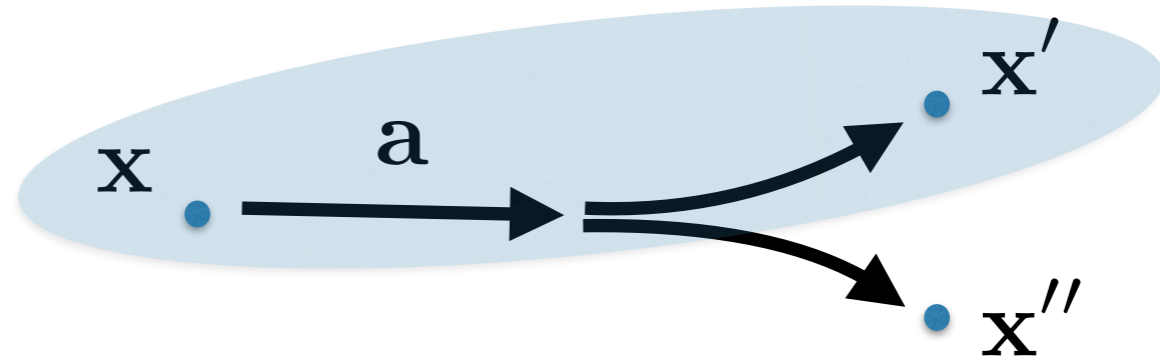
Both tasks formalised as reinforcement learning problems

States: $\mathbf{x} \in \mathcal{R}^n$

Actions: $\mathbf{a} \in \mathcal{R}^m$

Model: $p(\mathbf{x}' | \mathbf{x}, \mathbf{a})$

Rewards: $r(\mathbf{x}, \mathbf{a}, \mathbf{x}') \in \mathcal{R}$



Both tasks formalised as reinforcement learning problems

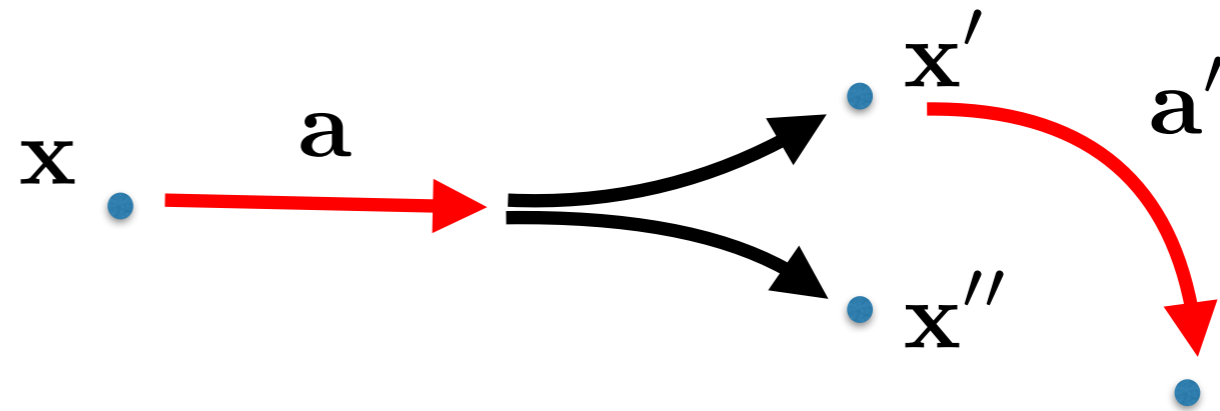
States: $\mathbf{x} \in \mathcal{R}^n$

Actions: $\mathbf{a} \in \mathcal{R}^m$

Model: $p(\mathbf{x}' | \mathbf{x}, \mathbf{a})$

Rewards: $r(\mathbf{x}, \mathbf{a}, \mathbf{x}') \in \mathcal{R}$

Policy: $\pi(\mathbf{a} | \mathbf{x})$



Both tasks formalised as reinforcement learning problems

States: $\mathbf{x} \in \mathcal{R}^n$

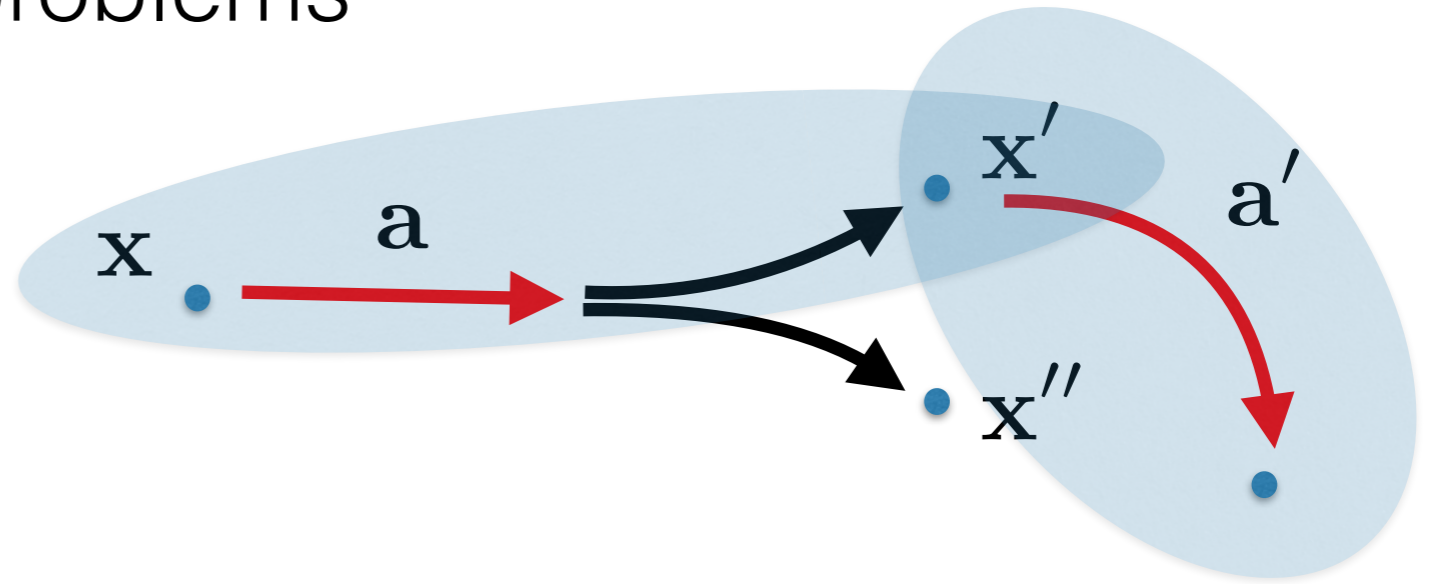
Actions: $\mathbf{a} \in \mathcal{R}^m$

Model: $p(\mathbf{x}' | \mathbf{x}, \mathbf{a})$

Rewards: $r(\mathbf{x}, \mathbf{a}, \mathbf{x}') \in \mathcal{R}$

Policy: $\pi(\mathbf{a} | \mathbf{x})$

Goal: $\pi^* = \arg \max_{\pi} J_{\pi}$ (e.g. $J_{\pi} = \mathbb{E} \left[\sum_{t=0}^T r_t \right]$)



Challenges in real tasks

States: $\mathbf{x} \in \mathcal{R}^n$ incomplete, noisy

Actions: $\mathbf{a} \in \mathcal{R}^m$ continuous high-dimensional

Model: $p(\mathbf{x}'|\mathbf{x}, \mathbf{a})$ inaccurate model

Rewards: $r(\mathbf{x}, \mathbf{a}, \mathbf{x}') \in \mathcal{R}$ hard to engineer

Policy: $\pi(\mathbf{a}|\mathbf{x})$ execution endanger the robot

Goal: $\pi^* = \arg \max_{\pi} J_{\pi}$ (e.g. $J_{\pi} = \mathbb{E} \left[\sum_{t=0}^T r_t \right]$)



Challenges in real tasks

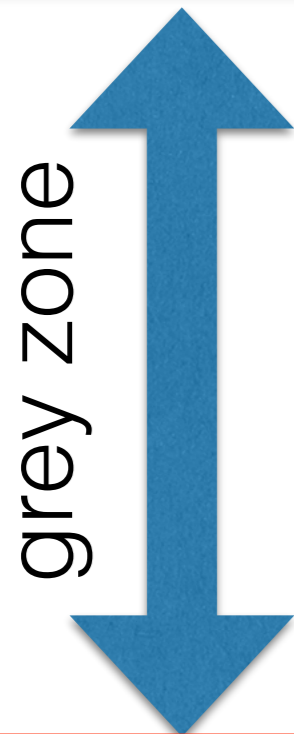
- Can I learn something without the model $p(\mathbf{x}'|\mathbf{x}, \mathbf{a})$ just from interactions?



Taxonomy of policy search methods

- Direct policy search (primal task)

e.g. gradient ascent for $\pi^* = \arg \max_{\pi} J_{\pi}$



Episodic REPS [Peters, 2010]

PILCO [Deisenroth, ICML 2011]

Actor-critic (e.g. DPG [Silver, JMLR 2014])

Deep Q-learning (e.g. [Mnih, Nature 2015])

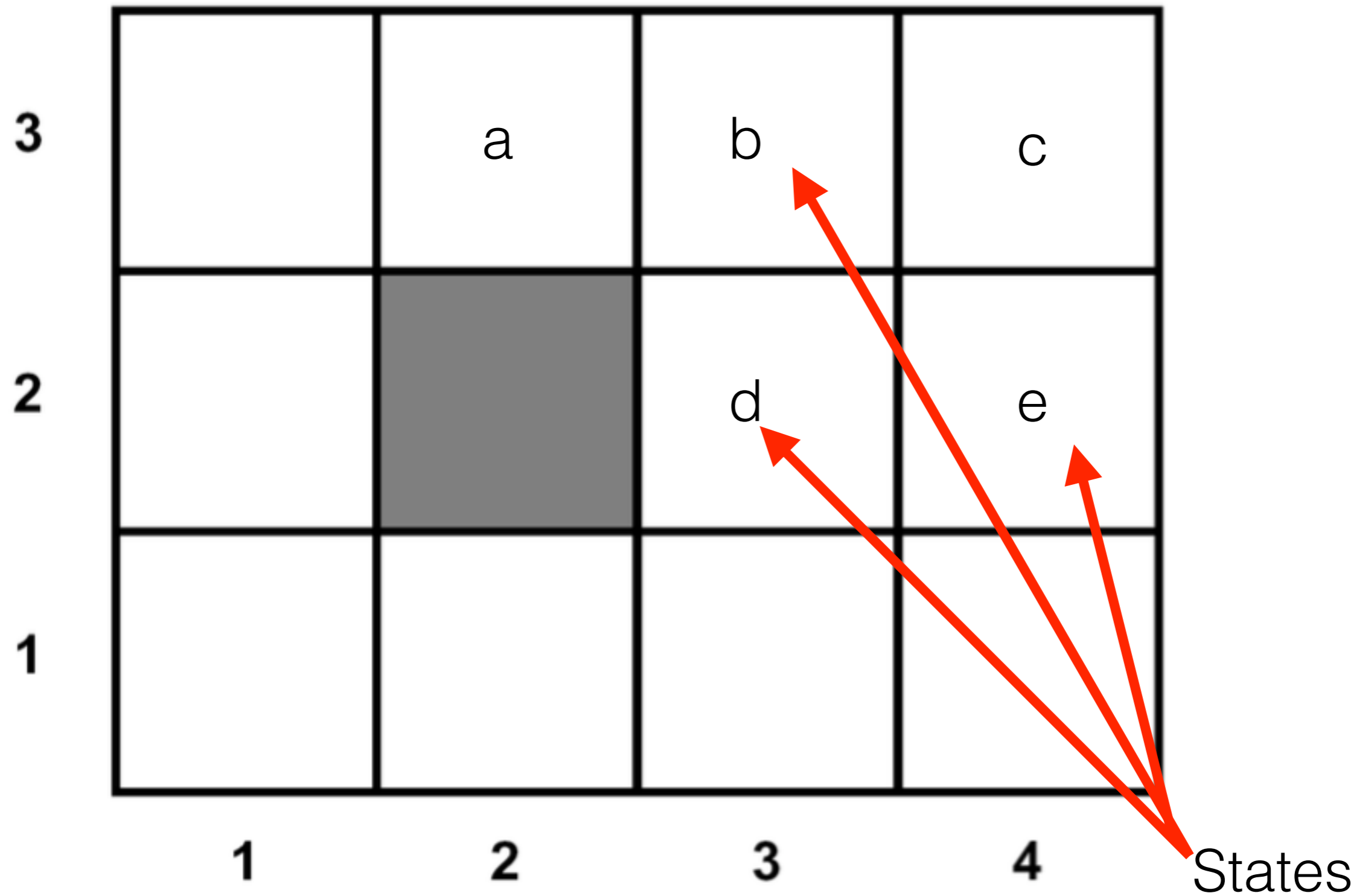
- Value-based methods (dual function [Kober, 2013])

e.g. search for $Q(\mathbf{x}, \mathbf{a}) = r(\mathbf{x}, \mathbf{a}, \mathbf{x}') + \gamma \max_{\mathbf{a}'} Q(\mathbf{x}', \mathbf{a}')$

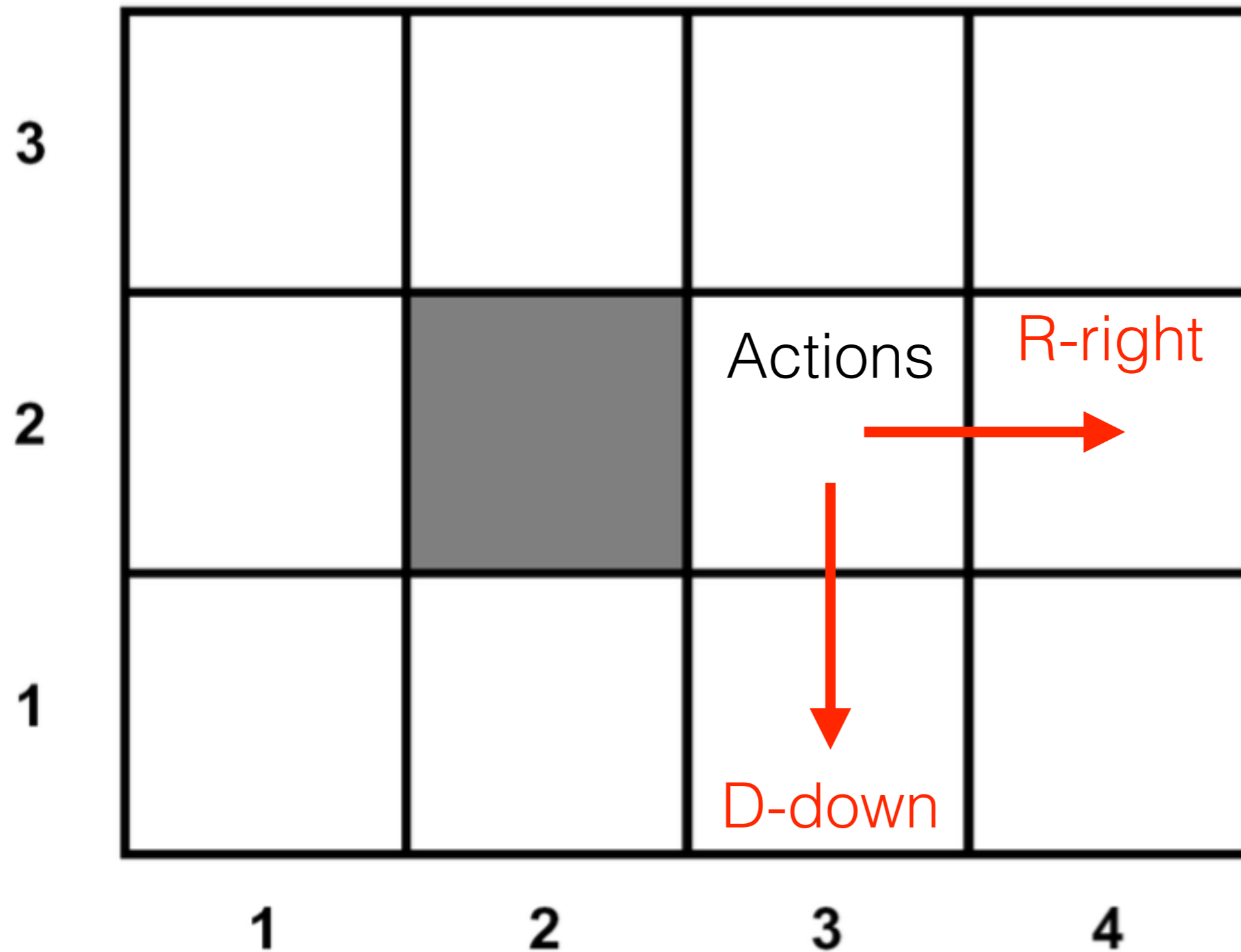
$\pi^* = \arg \max_a Q(\mathbf{x}, \mathbf{a})$



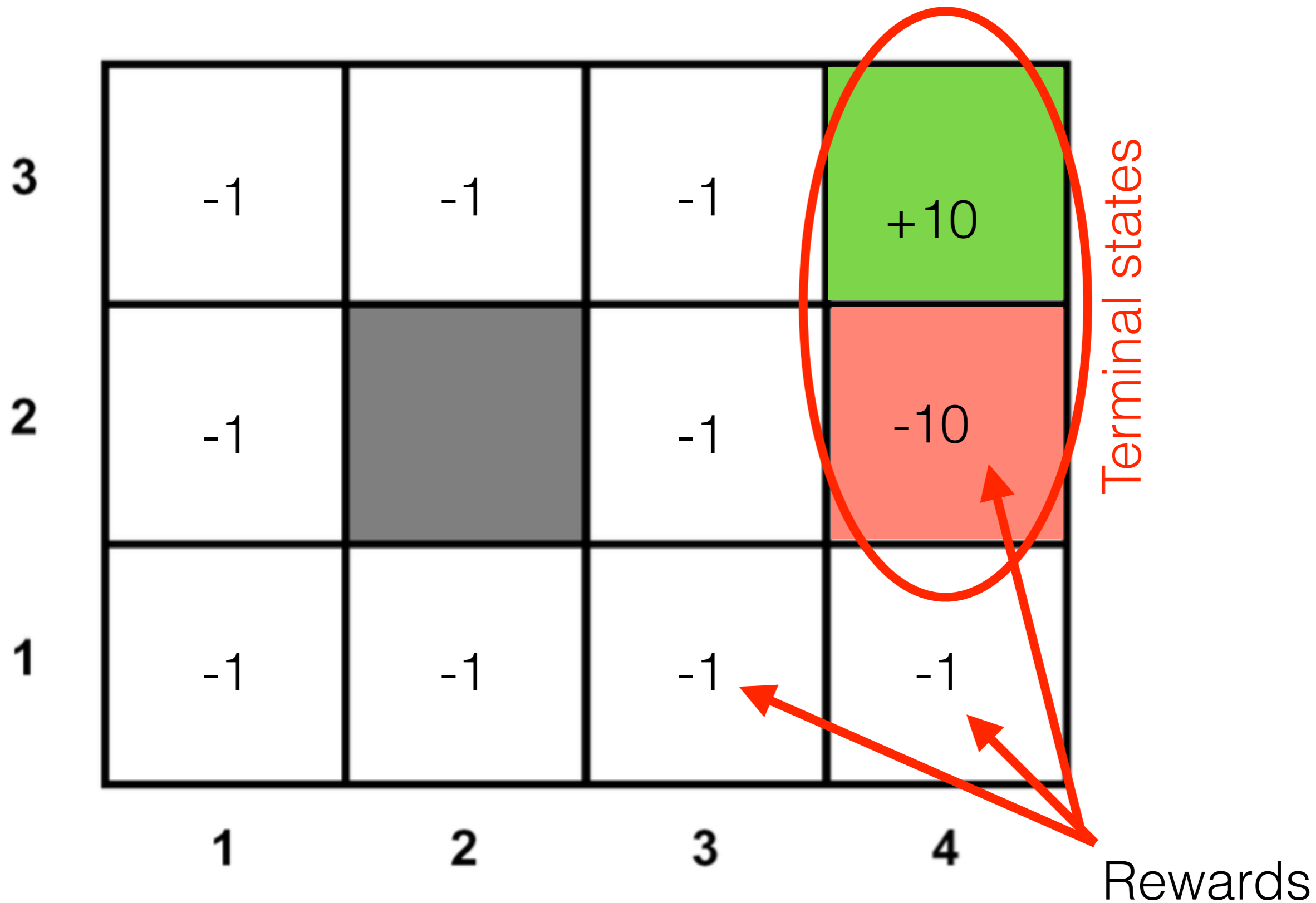
Value-based methods: Q-learning



Value-based methods: Q-learning



Value-based methods: Q-learning



	a	b	c
		d	e

State-action value function

$$Q(\mathbf{x}, \mathbf{u}) : X \times U \rightarrow \mathbb{R}$$

The best sum of rewards I can get, when following action u in state x and then controlling optimally

- Search for the Q , which satisfies Bellman equation

$$Q(\mathbf{x}, \mathbf{u}) = r(\mathbf{x}, \mathbf{u}, \mathbf{x}') + \max_{\mathbf{u}'} Q(\mathbf{x}', \mathbf{u}')$$



	a	b	c
		d	e

State-action value function

$$Q(\mathbf{x}, \mathbf{u}) : X \times U \rightarrow \mathbb{R}$$

The best sum of rewards I can get, when following action u in state x and then controlling optimally

- Search for the Q , which satisfies Bellman equation
- Once we find it, we can control optimally as follows:

$$Q(\mathbf{x}, \mathbf{u}) = r(\mathbf{x}, \mathbf{u}, \mathbf{x}') + \max_{\mathbf{u}'} Q(\mathbf{x}', \mathbf{u}')$$

$$\pi^*(\mathbf{x}) = \arg \max_{\mathbf{u}} Q(\mathbf{x}, \mathbf{u}) = \arg \max_{\pi} J_{\pi}$$



	a	b	c
		d	e

State-action value function

$$Q(\mathbf{x}, \mathbf{u}) : X \times U \rightarrow \mathbb{R}$$

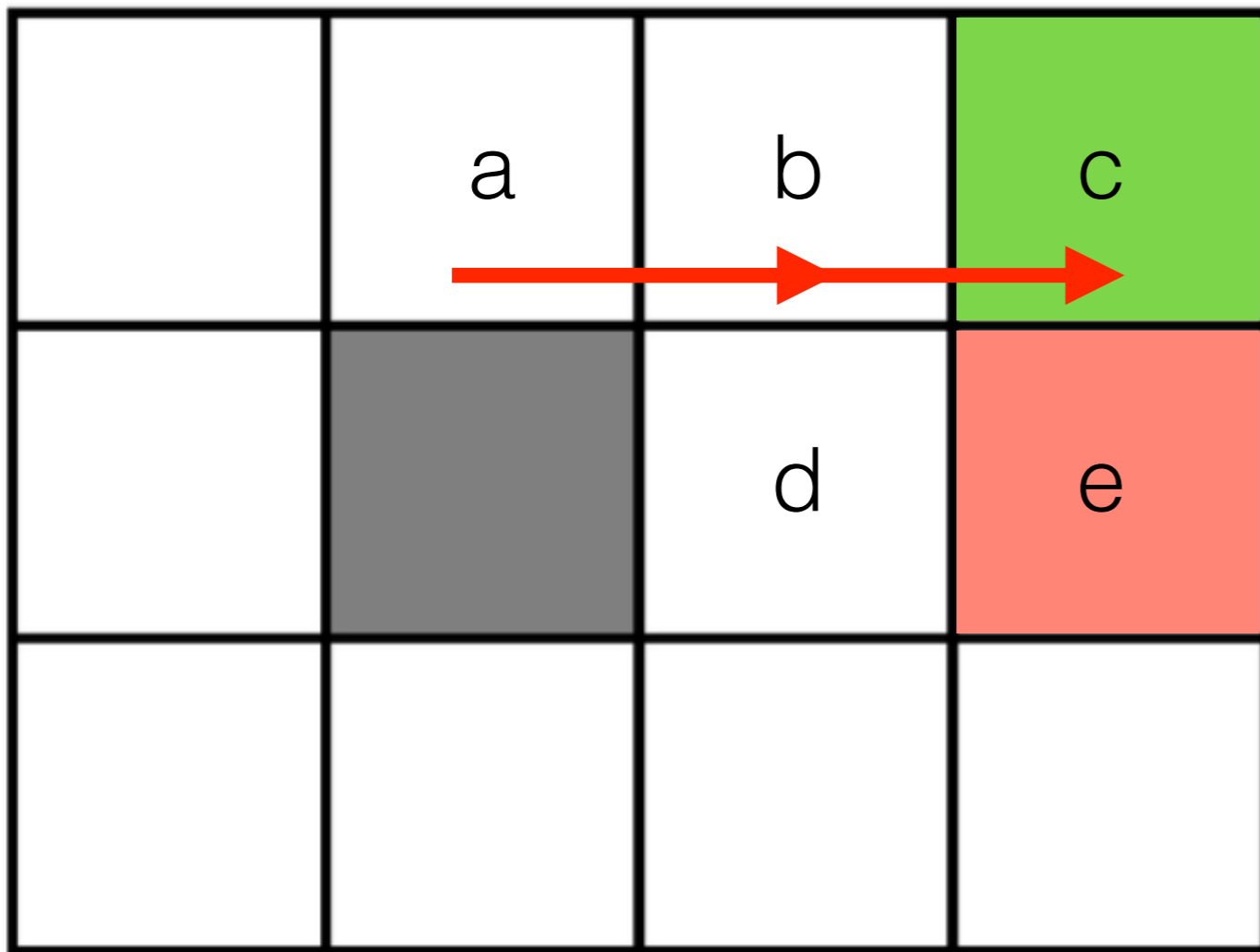
The best sum of rewards I can get, when following action u in state x and then controlling optimally

- Search for the Q , which satisfies Bellman equation

$$Q(\mathbf{x}, \mathbf{u}) = r(\mathbf{x}, \mathbf{u}, \mathbf{x}') + \max_{\mathbf{u}'} Q(\mathbf{x}', \mathbf{u}')$$
- Once we find it, we can control optimally as follows:

$$\pi^*(\mathbf{x}) = \arg \max_{\mathbf{u}} Q(\mathbf{x}, \mathbf{u}) = \arg \max_{\pi} J_{\pi}$$
- Search without model is based on collecting trajectories

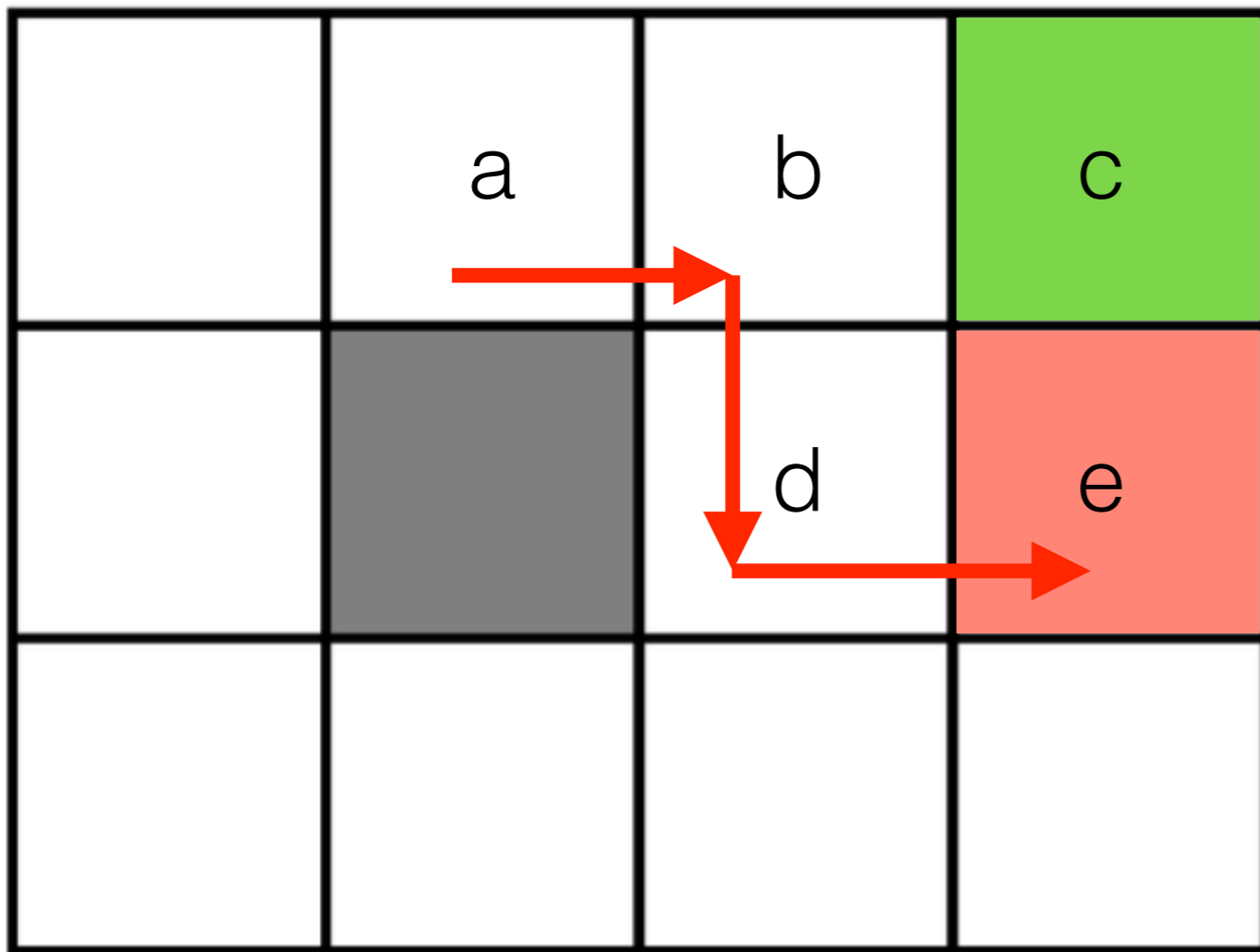




$\tau_1 :$
 $(a, R, -1), (b, R, -1), (c, R, 10)$

Q	R - right	D - down
a	?	?
b	?	?
c	?	?
d	?	?
e	?	?





$\tau_2 :$
 $(a, R, -1), (b, D, -1),$
 $(d, R, -1), (e, R, -10)$

Q	R - right	D - down
a	?	?
b	?	?
c	?	?
d	?	?
e	?	?

Having a trajectory, each transition gives one equation



	a	b	c
		d	e

$\tau_2 :$

$(a, R, -1), (b, D, -1),$

$(d, R, -1), (e, R, -10)$



$$Q(b, D) = r(b) + \max_{\mathbf{u}} Q(d, \mathbf{u})$$

Q	R - right	D - down
a	?	?
b	?	?
c	?	?
d	?	?
e	?	?

Having a trajectory, each transition gives one equation



	a	b	c
		d	e

$\tau_2 :$

$(a, R, -1), (b, D, -1),$
 $(d, R, -1), (e, R, -10)$

$$Q(b, D) = r(b) + \max_{\mathbf{u}} Q(d, \mathbf{u})$$

$$Q(d, R) = r(d) + \max_{\mathbf{u}} Q(e, \mathbf{u})$$

$$Q(a, R) = r(a) + \max_{\mathbf{u}} Q(b, \mathbf{u})$$

Q	R - right	D - down
a	?	?
b	?	?
c	?	?
d	?	?
e	?	?

Having a trajectory, each transition gives one equation



	a	b	c
		d	e

$\tau_2 :$

$(a, R, -1), (b, D, -1),$
 $(d, R, -1), (e, R, -10)$

$Q(e, R) = r(e)$

$Q(b, D) = r(b) + \max_{\mathbf{u}} Q(d, \mathbf{u})$

$Q(d, R) = r(d) + \max_{\mathbf{u}} Q(e, \mathbf{u})$

$Q(a, R) = r(a) + \max_{\mathbf{u}} Q(b, \mathbf{u})$

Q	R - right	D - down
a	?	?
b	?	?
c	?	?
d	?	?
e	?	?

Having a trajectory, each transition gives one equation



	a	b	c
		d	e

$\tau_2 :$

$(a, R, -1), (b, D, -1),$
 $(d, R, -1), (e, R, -10)$

$$Q(e, R) = r(e)$$

$$Q(b, D) = r(b) + \max_{\mathbf{u}} Q(d, \mathbf{u})$$

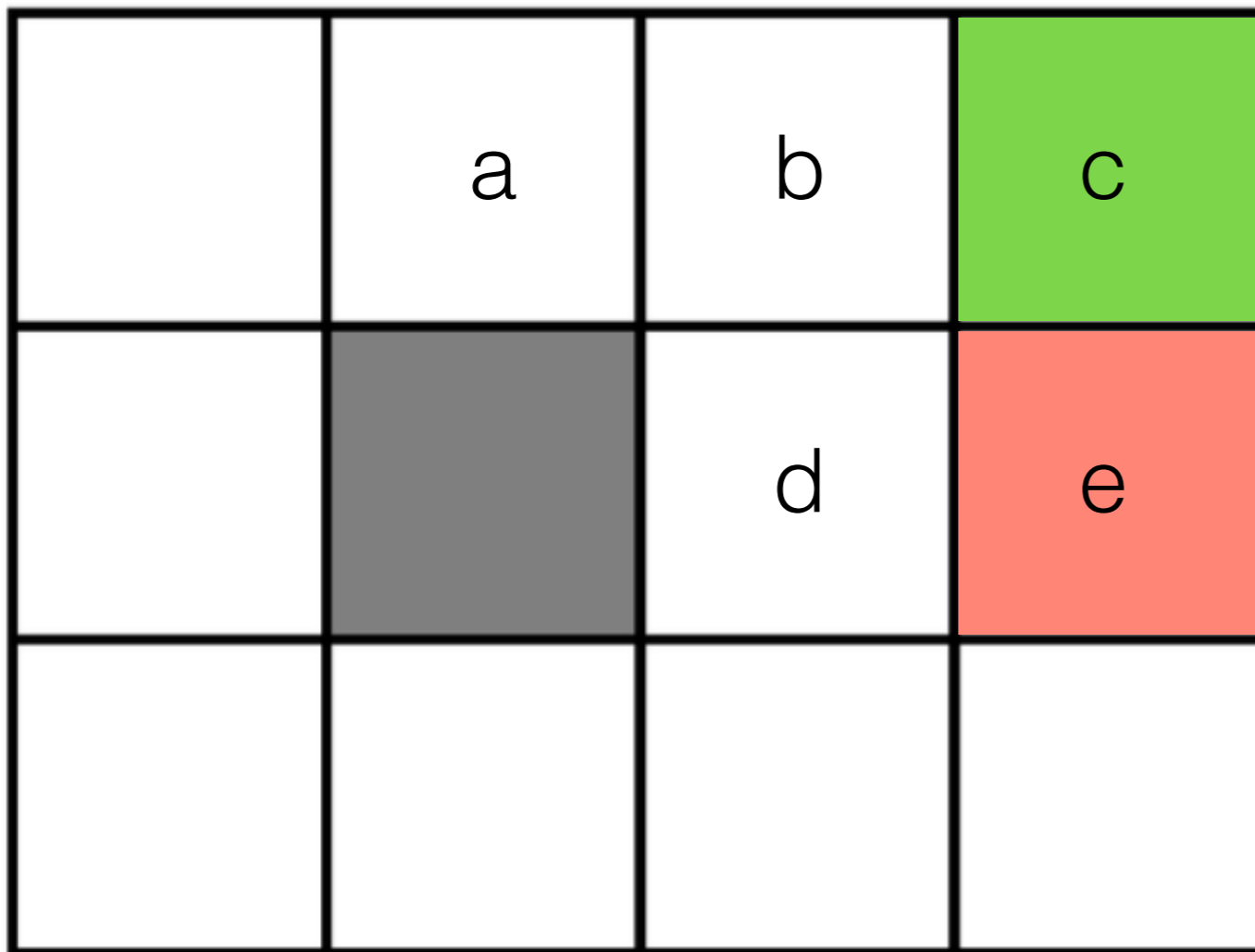
$$Q(d, R) = r(d) + \max_{\mathbf{u}} Q(e, \mathbf{u})$$

$$Q(a, R) = r(a) + \max_{\mathbf{u}} Q(b, \mathbf{u})$$

Q	R - right	D - down
a	?	?
b	?	?
c	?	?
d	?	?
e	?	?

Having a trajectory, each transition gives one equation





$\tau_2 :$
 $(a, R, -1), (b, D, -1),$
 $(d, R, -1), (e, R, -10)$

$$Q(e, R) = r(e)$$

$$Q(b, D) = r(b) + \max_{\mathbf{u}} Q(d, \mathbf{u})$$

$$Q(d, R) = r(d) + \max_{\mathbf{u}} Q(e, \mathbf{u})$$

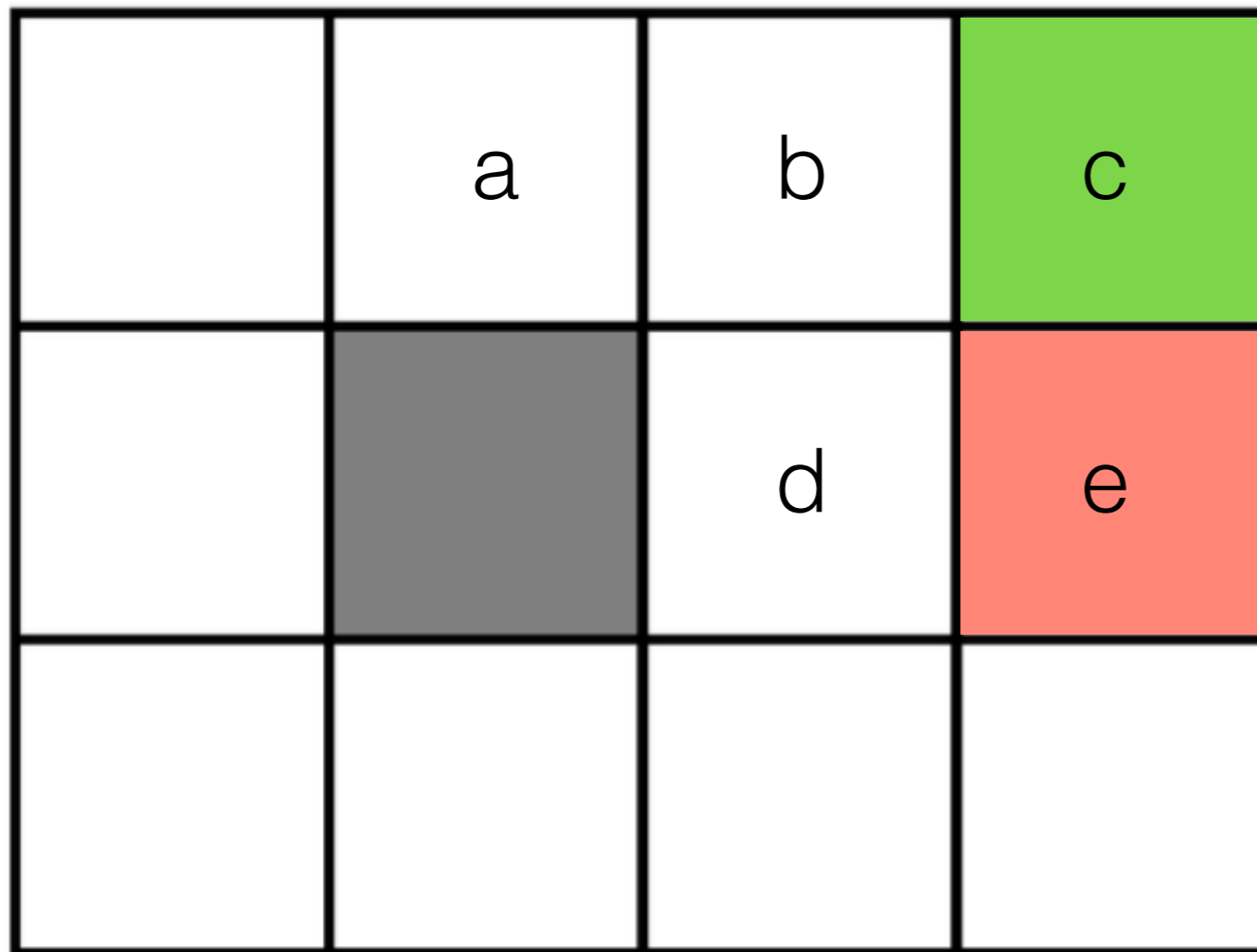
$$Q(a, R) = r(a) + \max_{\mathbf{u}} Q(b, \mathbf{u})$$

Q	R - right	D - down
a	?	?
b	?	?
c	?	?
d	?	?
e	?	?

 unknowns

Having a trajectory, each transition gives one equation





$\tau_2 :$
 $(a, R, -1), (b, D, -1),$
 $(d, R, -1), (e, R, -10)$

$$Q(e, R) = r(e)$$

$$Q(b, D) = r(b) + \max_{\mathbf{u}} Q(d, \mathbf{u})$$

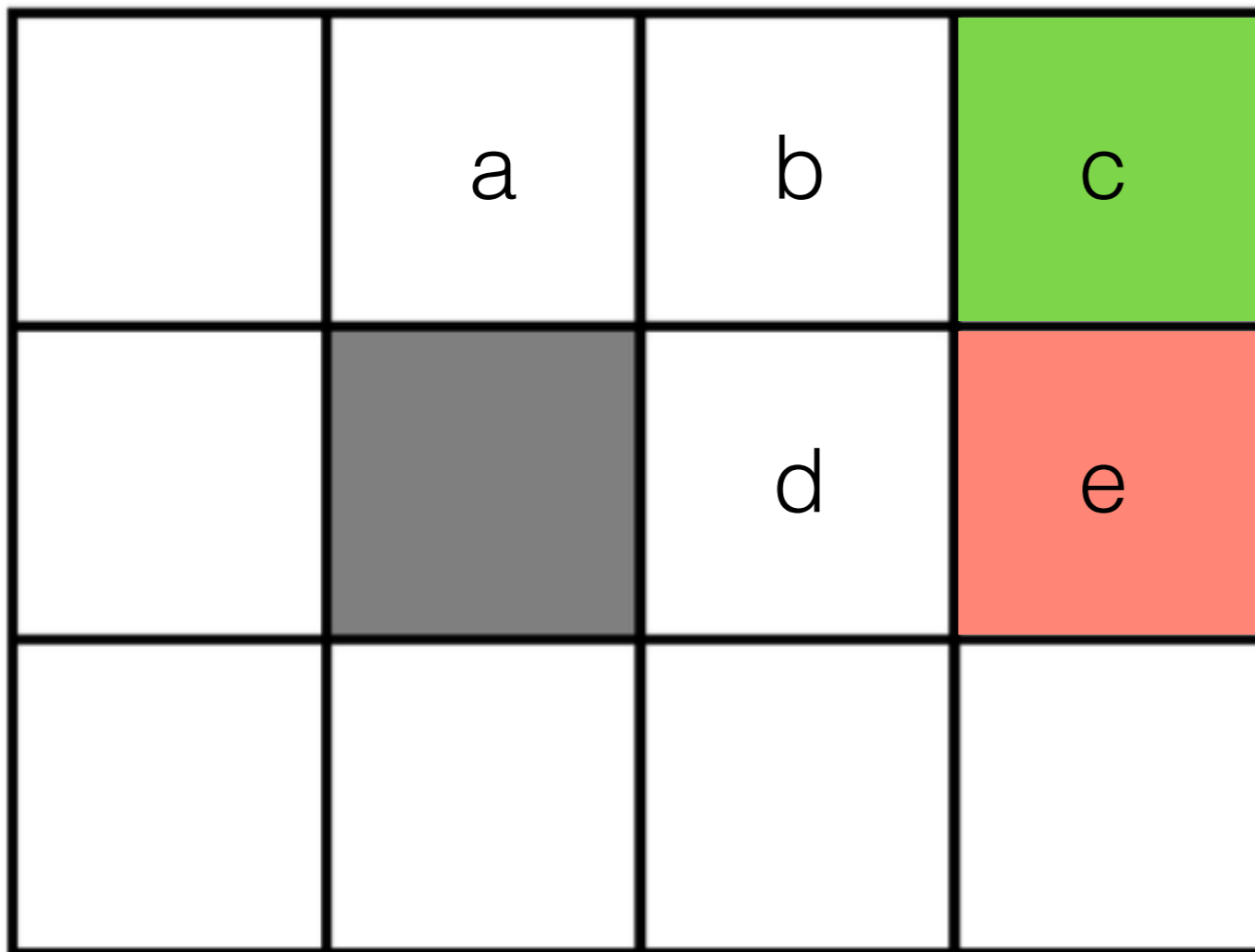
$$Q(d, R) = r(d) + \max_{\mathbf{u}} Q(e, \mathbf{u})$$

$$Q(a, R) = r(a) + \max_{\mathbf{u}} Q(b, \mathbf{u})$$

Q	R - right	D - down
a	0	0
b	0	0
c	0	0
d	0	0
e	0	0

(1) Substitute transitions and current Q-values to the right side and solve for left side.





$\tau_2 :$

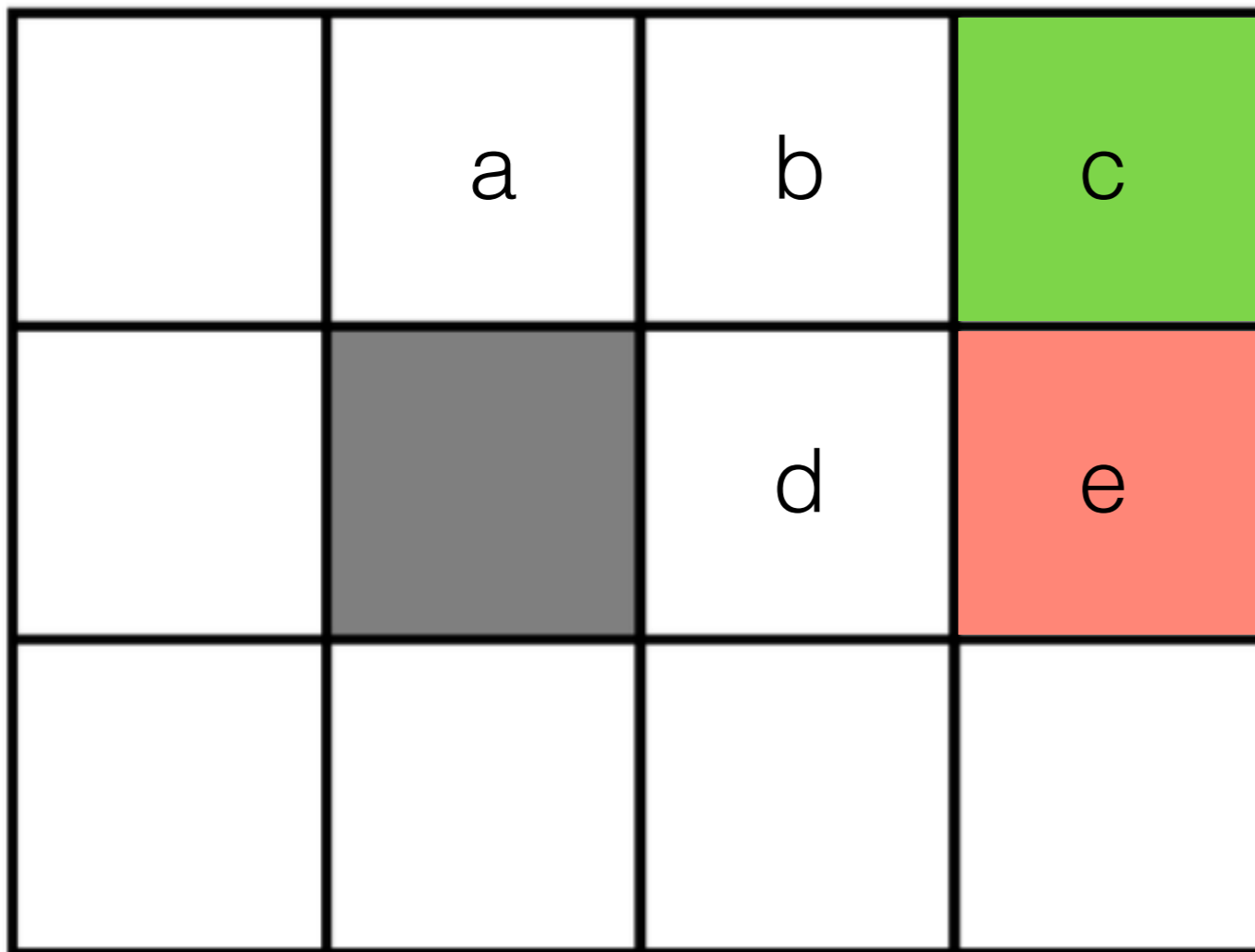
$(a, R, -1), (b, D, -1),$
 $(d, R, -1), (e, R, -10)$

$Q(e, R)$	$=$	-10
$Q(b, D)$	$=$	$r(b) + \max_{\mathbf{u}} Q(d, \mathbf{u})$
$Q(d, R)$	$=$	$r(d) + \max_{\mathbf{u}} Q(e, \mathbf{u})$
$Q(a, R)$	$=$	$r(a) + \max_{\mathbf{u}} Q(b, \mathbf{u})$

Q	R - right	D - down
a	0	0
b	0	0
c	0	0
d	0	0
e	-10	0

(1) Substitute transitions and current Q-values to the right side and solve for left side.





$\tau_2 :$

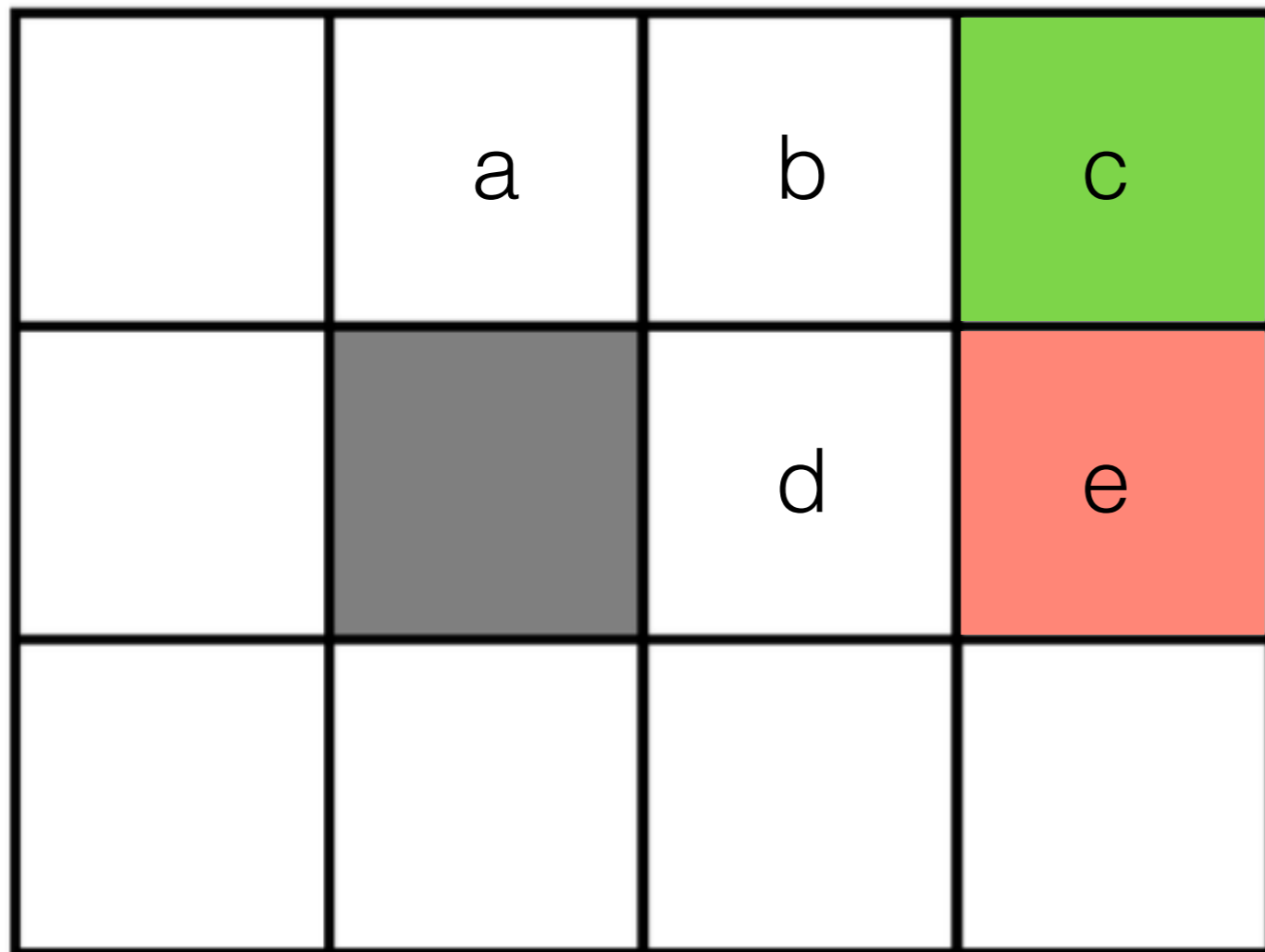
$(a, R, -1), (b, D, -1),$
 $(d, R, -1), (e, R, -10)$

$Q(e, R)$	=	-10
$Q(b, D)$	=	???
$Q(d, R)$	=	$r(d) + \max_{\mathbf{u}} Q(e, \mathbf{u})$
$Q(a, R)$	=	$r(a) + \max_{\mathbf{u}} Q(b, \mathbf{u})$

Q	R - right	D - down
a	0	0
b	0	-1
c	0	0
d	0	0
e	-10	0

(1) Substitute transitions and current Q-values to the right side and solve for left side.





$\tau_2 :$

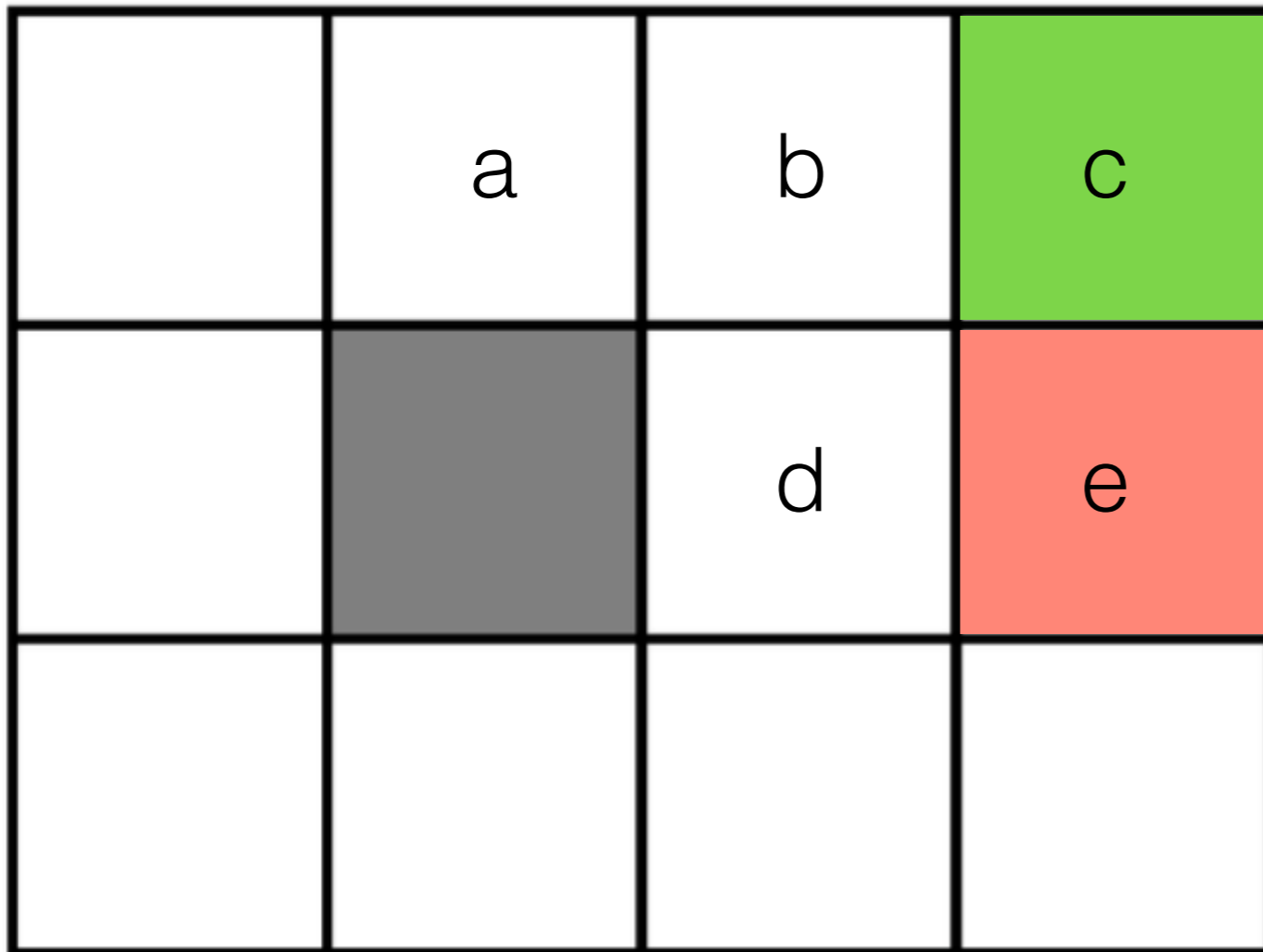
$(a, R, -1), (b, D, -1),$
 $(d, R, -1), (e, R, -10)$

$Q(e, R)$	=	-10
$Q(b, D)$	=	-1
$Q(d, R)$	=	$r(d) + \max_{\mathbf{u}} Q(e, \mathbf{u})$
$Q(a, R)$	=	$r(a) + \max_{\mathbf{u}} Q(b, \mathbf{u})$

Q	R - right	D - down
a	0	0
b	0	-1
c	0	0
d	0	0
e	-10	0

(1) Substitute transitions and current Q-values to the right side and solve for left side.





$\tau_2 :$

$(a, R, -1), (b, D, -1),$

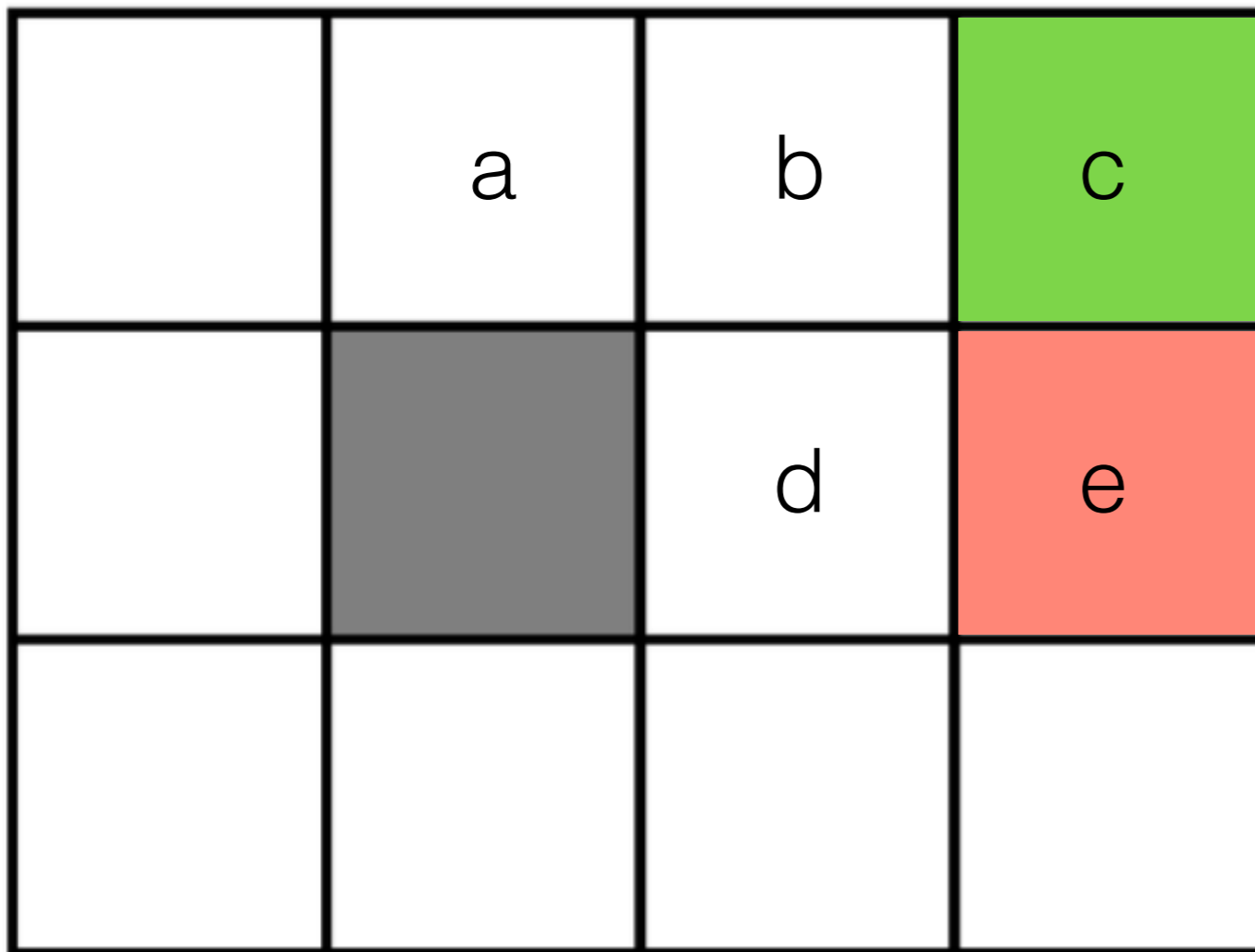
$(d, R, -1), (e, R, -10)$

$Q(e, R)$	=	-10
$Q(b, D)$	=	-1
$Q(d, R)$	=	-1
$Q(a, R)$	=	$r(a) + \max_{\mathbf{u}} Q(b, \mathbf{u})$

Q	R - right	D - down
a	0	0
b	0	-1
c	0	0
d	-1	0
e	-10	0

(1) Substitute transitions and current Q-values to the right side and solve for left side.





$\tau_2 :$

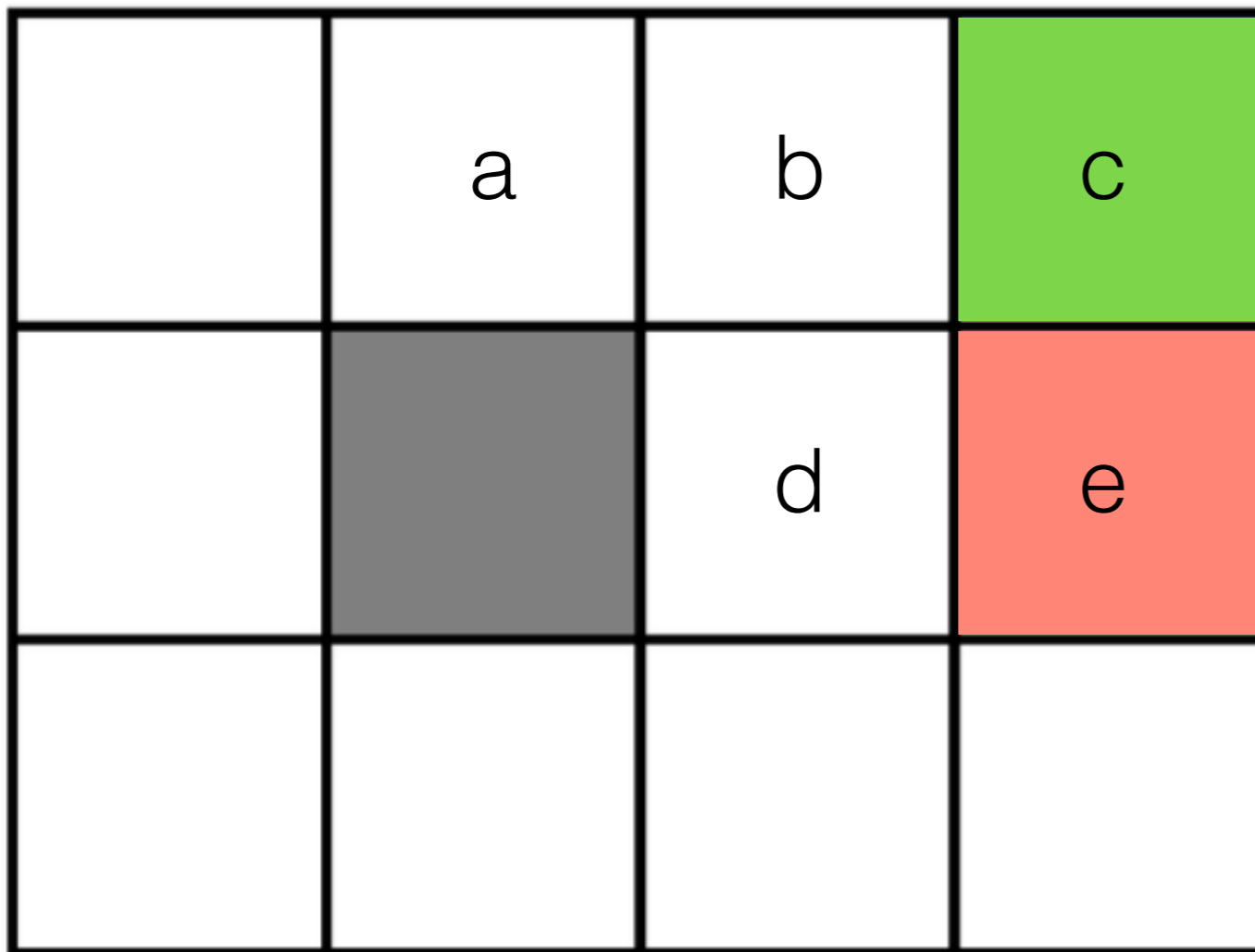
$(a, R, -1), (b, D, -1),$
 $(d, R, -1), (e, R, -10)$

$Q(e, R)$	=	-10
$Q(b, D)$	=	-1
$Q(d, R)$	=	-1
$Q(a, R)$	=	-1

Q	R - right	D - down
a	-1	0
b	0	-1
c	0	0
d	-1	0
e	-10	0

(1) Substitute transitions and current Q-values to the right side and solve for left side.





$\tau_2 :$

$(a, R, -1), (b, D, -1),$
 $(d, R, -1), (e, R, -10)$

$$Q(e, R) = r(e)$$

$$Q(b, D) = r(b) + \max_{\mathbf{u}} Q(d, \mathbf{u})$$

$$Q(d, R) = r(d) + \max_{\mathbf{u}} Q(e, \mathbf{u})$$

$$Q(a, R) = r(a) + \max_{\mathbf{u}} Q(b, \mathbf{u})$$

Q	R - right	D - down
a	-1	0
b	0	-1
c	0	0
d	-1	0
e	-10	0

- (1) Substitute transitions and current Q-values to the right side and solve for left side.
 (2) Repeat several times



	a	b	c
		d	e

$\tau_2 :$

$(a, R, -1), (b, D, -1),$
 $(d, R, -1), (e, R, -10)$

$$Q(e, R) = r(e)$$

$$Q(b, R) = r(b) + \max_{\mathbf{u}} Q(d, \mathbf{u})$$

$$Q(d, R) = r(d) + \max_{\mathbf{u}} Q(e, \mathbf{u})$$

$$Q(a, R) = r(a) + \max_{\mathbf{u}} Q(b, \mathbf{u})$$

Q	R - right	D - down
a	-1	0
b	0	-1
c	0	0
d	-1	0
e	-10	0

(1) Substitute transitions and current Q-values to the right side and solve for left side.
 (2) Repeat several times (search for the fixed point of the Bellman operator)

$$Q = \mathcal{B}(Q)$$



	a	b	c
		d	e

$\tau_2 :$

$(a, R, -1), (b, D, -1),$
 $(d, R, -1), (e, R, -10)$

$$Q(e, R) = r(e)$$

$$Q(b, R) = r(b) + \max_{\mathbf{u}} Q(d, \mathbf{u})$$

$$Q(d, R) = r(d) + \max_{\mathbf{u}} Q(e, \mathbf{u})$$

$$Q(a, R) = r(a) + \max_{\mathbf{u}} Q(b, \mathbf{u})$$

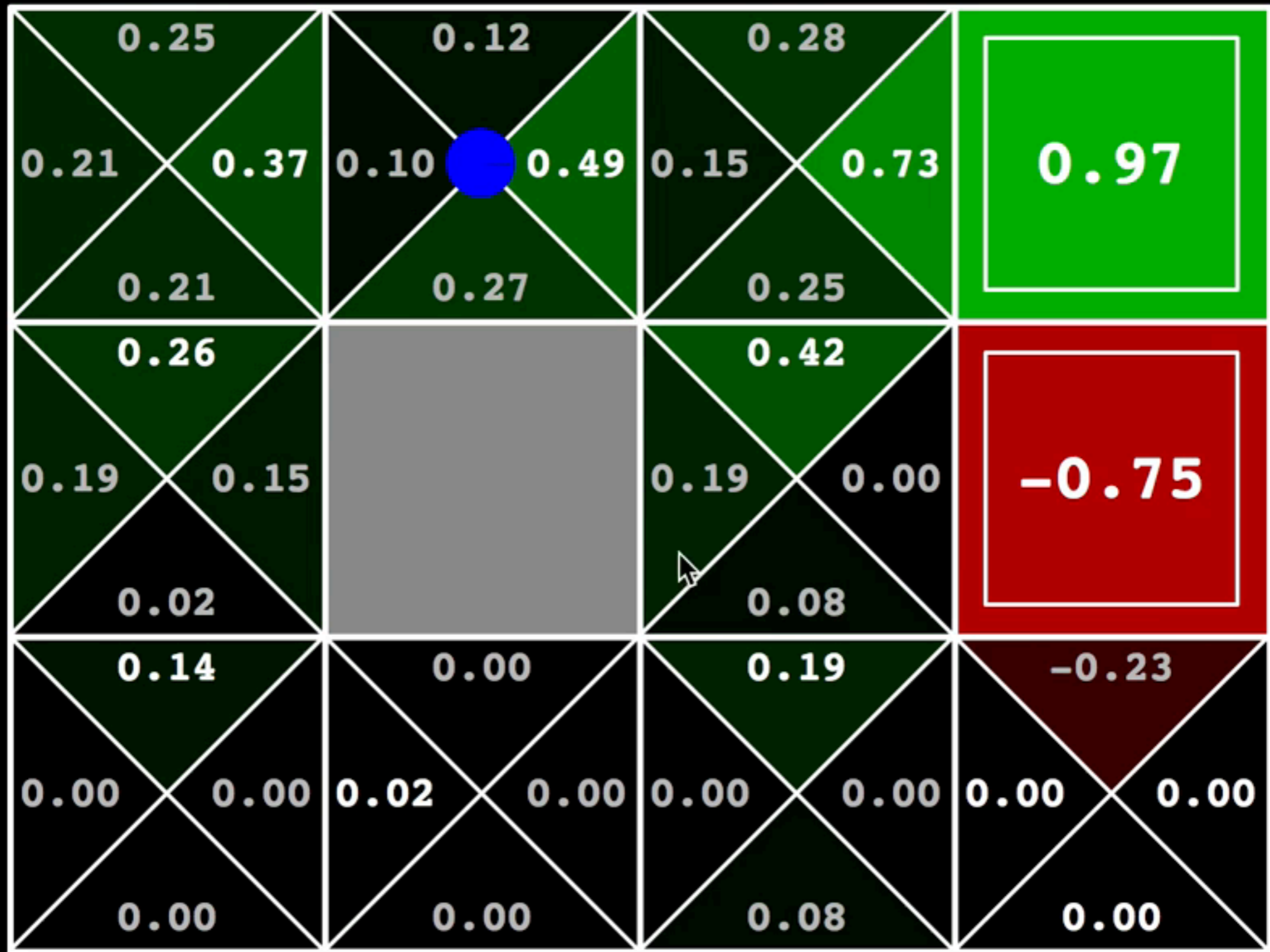
Iterations of the Bellman operator converge to a fixed point !!!

(1) Substitute transitions and current Q-values to the right side and solve for left side.

(2) Repeat several times (search for the fixed point of the Bellman operator)

$$Q = \mathcal{B}(Q)$$





CURRENT Q-VALUES



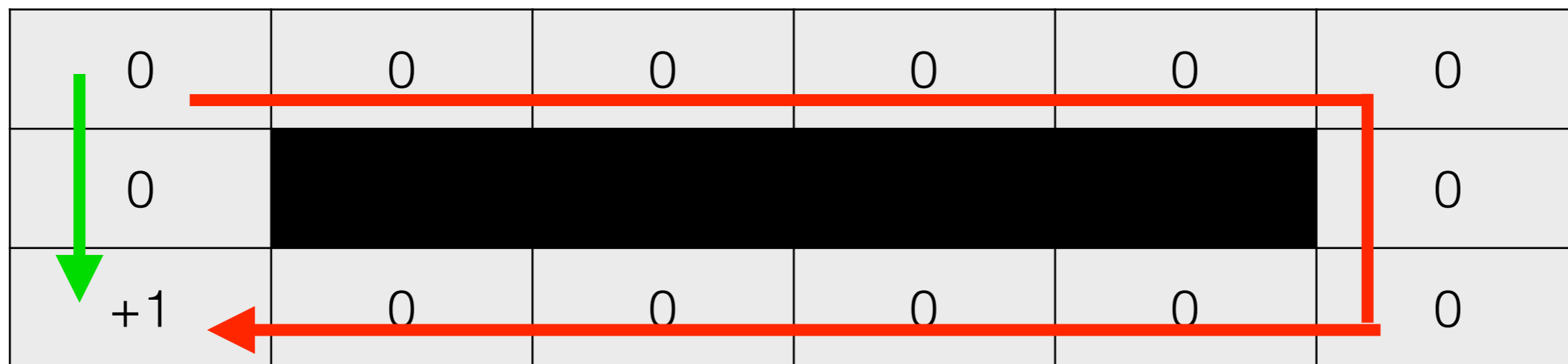
Bellman equation

$$Q(\mathbf{x}, \mathbf{u}) = r(\mathbf{x}, \mathbf{u}, \mathbf{x}') + \max_{\mathbf{u}'} Q(\mathbf{x}', \mathbf{u}')$$

reward for transition

the best you can do from
the following state

Which path is better?



Bellman equation

$$Q(\mathbf{x}, \mathbf{u}) = r(\mathbf{x}, \mathbf{u}, \mathbf{x}') + \gamma \max_{\mathbf{u}'} Q(\mathbf{x}', \mathbf{u}')$$

reward for transition

the best you can do from
the following state

discount factor $\gamma \in [0; 1]$

0	0	0	0	0	0	0
0						0
+1	0	0	0	0	0	0



Q-learning

1. Collect trajectories $\tau_1, \tau_2, \tau_3, \dots$
2. Solve $Q(\mathbf{x}, \mathbf{u}) = r(\mathbf{x}, \mathbf{u}, \mathbf{x}') + \gamma \max_{\mathbf{u}'} Q(\mathbf{x}', \mathbf{u}')$
3. Repeat from 1



Q-learning

1. Collect trajectories $\tau_1, \tau_2, \tau_3, \dots$
 2. Solve $Q(\mathbf{x}, \mathbf{u}) = r(\mathbf{x}, \mathbf{u}, \mathbf{x}') + \gamma \max_{\mathbf{u}'} Q(\mathbf{x}', \mathbf{u}')$
 3. Repeat from 1
- Curse of dimensionality



Q-learning

1. Collect trajectories $\tau_1, \tau_2, \tau_3, \dots$
 2. Solve $Q(\mathbf{x}, \mathbf{u}) = r(\mathbf{x}, \mathbf{u}, \mathbf{x}') + \gamma \max_{\mathbf{u}'} Q(\mathbf{x}', \mathbf{u}')$
 3. Repeat from 1
- Curse of dimensionality
 - Replace table $Q(\mathbf{x}, \mathbf{u})$ by function $Q_\theta(\mathbf{x}, \mathbf{u})$



Q-learning

1. Collect trajectories $\tau_1, \tau_2, \tau_3, \dots$
 2. Solve $Q(\mathbf{x}, \mathbf{u}) = r(\mathbf{x}, \mathbf{u}, \mathbf{x}') + \gamma \max_{\mathbf{u}'} Q(\mathbf{x}', \mathbf{u}')$
 3. Repeat from 1
- Curse of dimensionality
 - Replace table $Q(\mathbf{x}, \mathbf{u})$ by function $Q_\theta(\mathbf{x}, \mathbf{u})$

Approximate Q-learning

1. Collect trajectories $\tau_1, \tau_2, \tau_3, \dots$, initialize $\theta = \text{rand}$
2. Estimate $\mathbf{y} = r(\mathbf{x}, \mathbf{u}, \mathbf{x}') + \gamma \max_{\mathbf{u}'} Q_\theta(\mathbf{x}', \mathbf{u}')$
3. Update parameters by learning

$$\arg \min_{\theta} \sum_{\mathbf{x}, \mathbf{u}, \mathbf{y}} \|Q_\theta(\mathbf{x}, \mathbf{u}) - \mathbf{y}\|$$

4. Repeat from 2
5. Repeat from 1



Q-learning

1. Collect trajectories $\tau_1, \tau_2, \tau_3, \dots$
 2. Solve $Q(\mathbf{x}, \mathbf{u}) = r(\mathbf{x}, \mathbf{u}, \mathbf{x}') + \gamma \max_{\mathbf{u}'} Q(\mathbf{x}', \mathbf{u}')$
 3. Repeat from 1
- Curse of dimensionality
 - Replace table $Q(\mathbf{x}, \mathbf{u})$ by function $Q_\theta(\mathbf{x}, \mathbf{u})$

Approximate Q-learning

1. Collect trajectories $\tau_1, \tau_2, \tau_3, \dots$, initialize $\theta = \text{rand}$
2. Estimate $\mathbf{y} = r(\mathbf{x}, \mathbf{u}, \mathbf{x}') + \gamma \max_{\mathbf{u}'} Q_\theta(\mathbf{x}', \mathbf{u}')$
3. Update parameters by learning

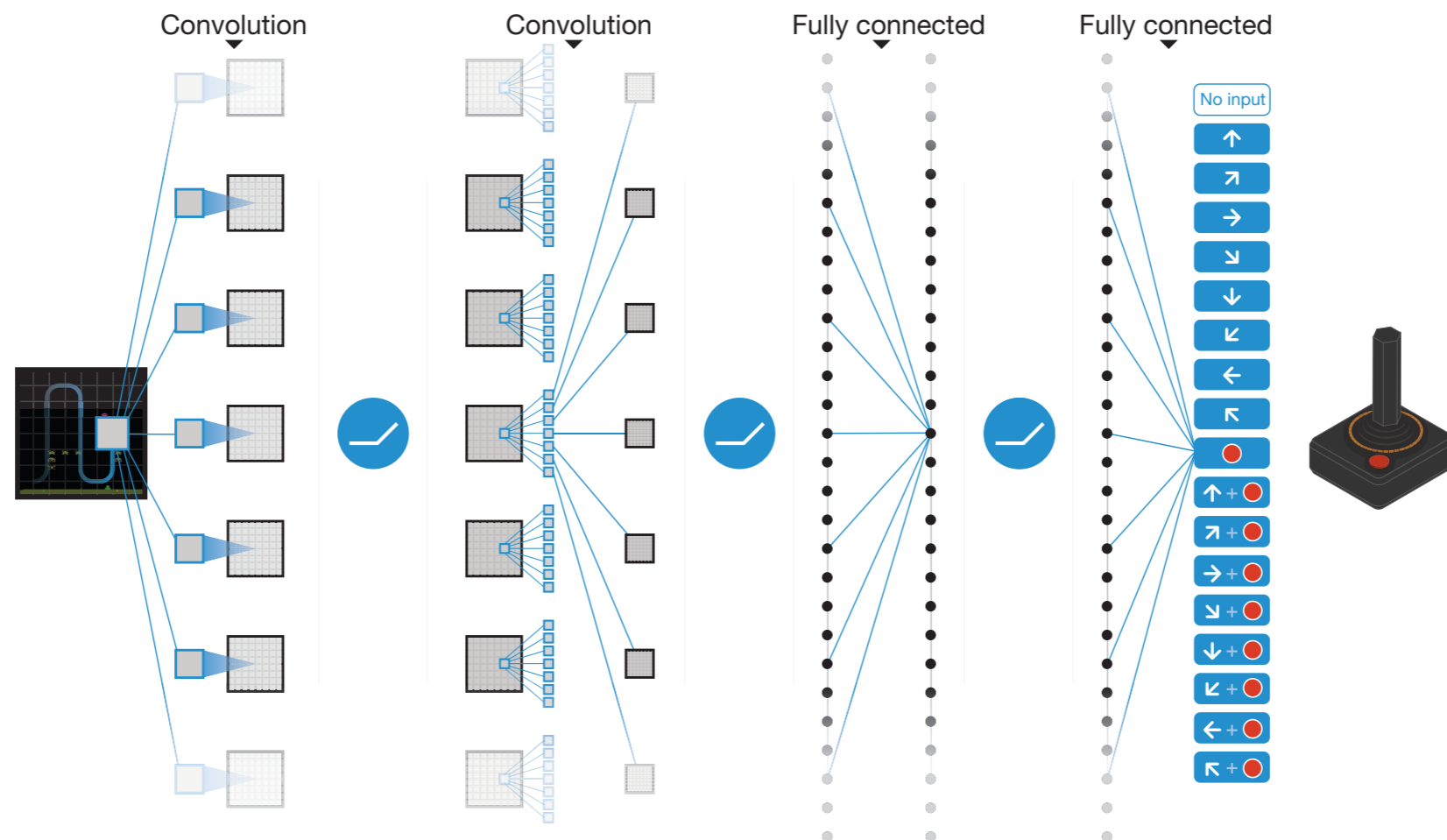
$$\arg \min_{\theta} \sum_{\mathbf{x}, \mathbf{u}, \mathbf{y}} \|Q_\theta(\mathbf{x}, \mathbf{u}) - \mathbf{y}\|$$

4. Repeat from 2
 5. Repeat from 1
- Approximated Q-learning does not have to converge to a fixed-point !!!**



Mnih et al. Nature 2015

- 2600 atari games
- **state space:** pixels (e.g. VGA resolution)
- **action space:** discrete joystic actions (8 direction + 8 direction with button + neutral action)
- replay buffer (decorrelates samples to be “more i.i.d”)
- two Q-networks (suppress oscillations)

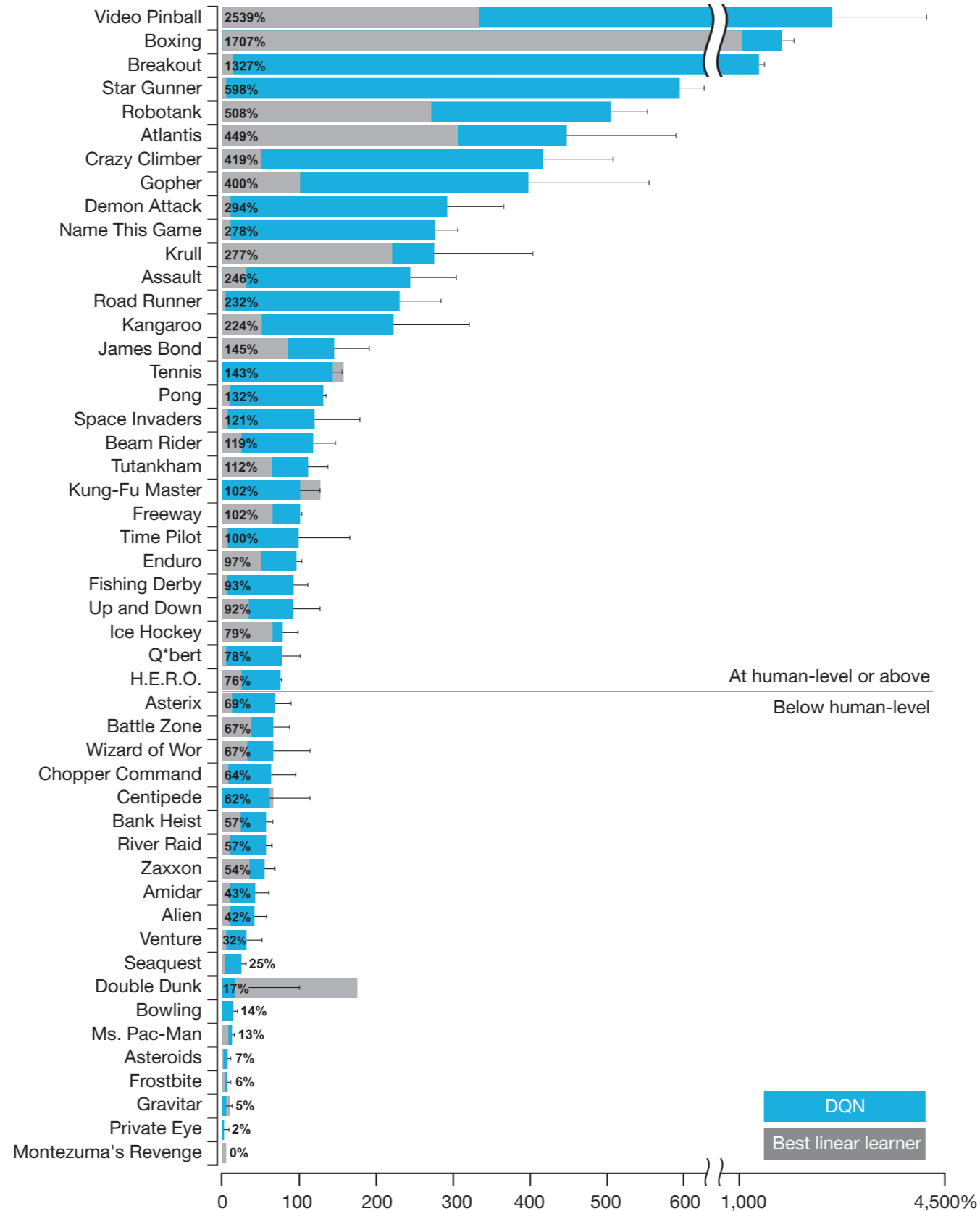


Mnih et al. Nature 2015

- 2600 atari games
- **state space:** pixels (e.g. VGA resolution)
- **action space:** discrete joystic actions (8 directions + 8 directions with button)
- collection of control tasks: <https://gym.openai.com>



Mnih et al. Nature 2015

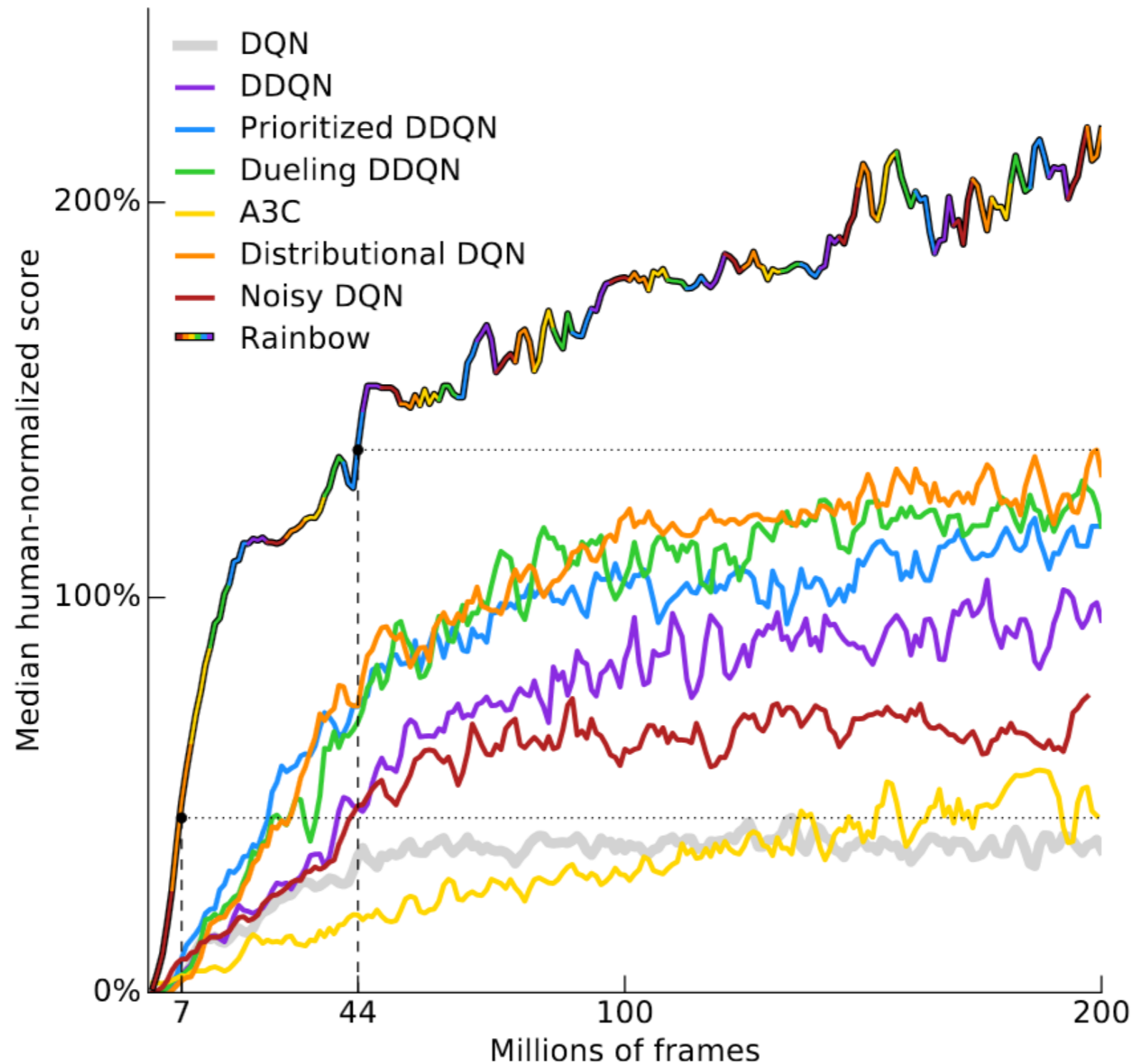


Czech Technical University in Prague

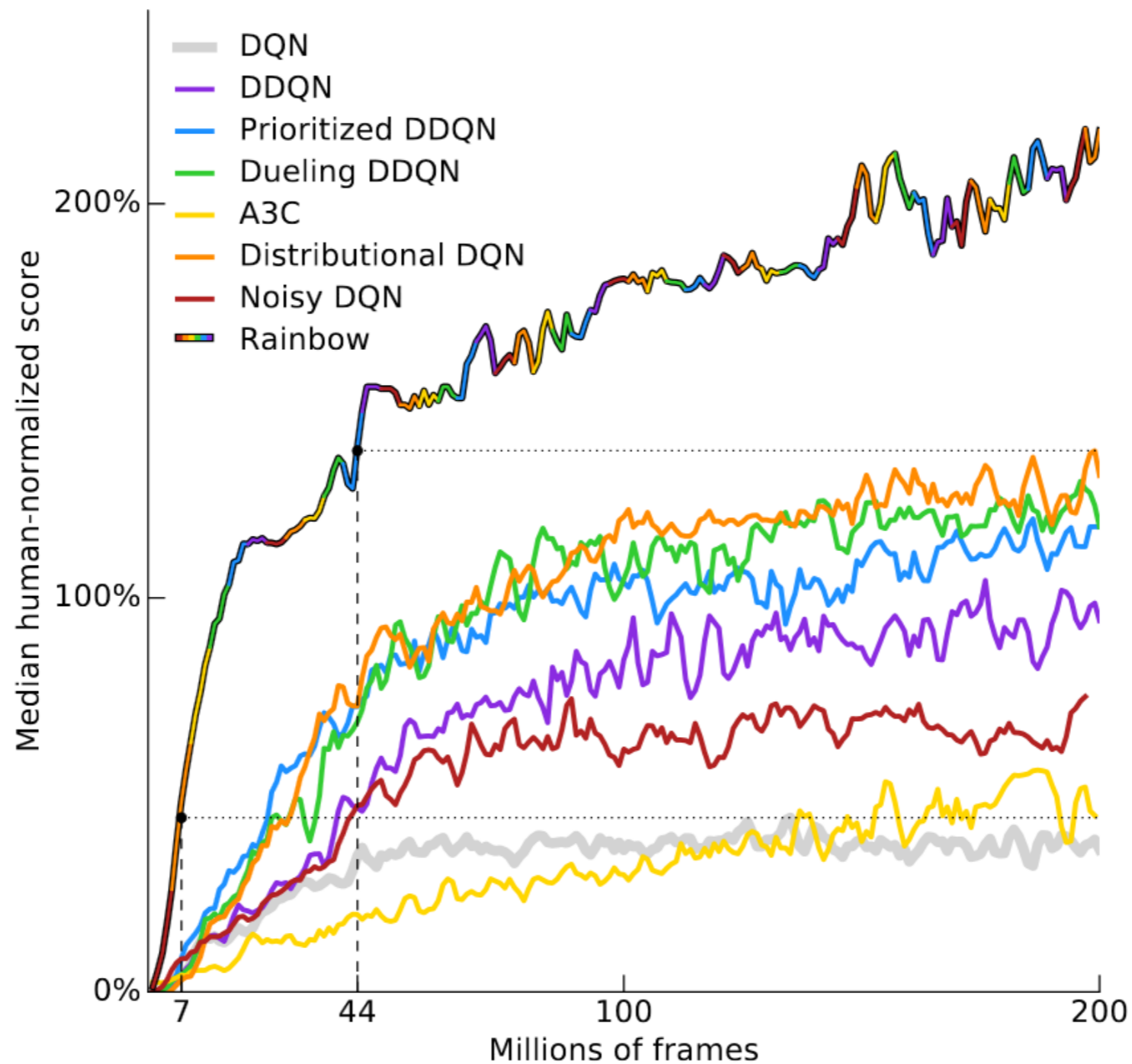
Faculty of Electrical Engineering, Department of Cybernetics



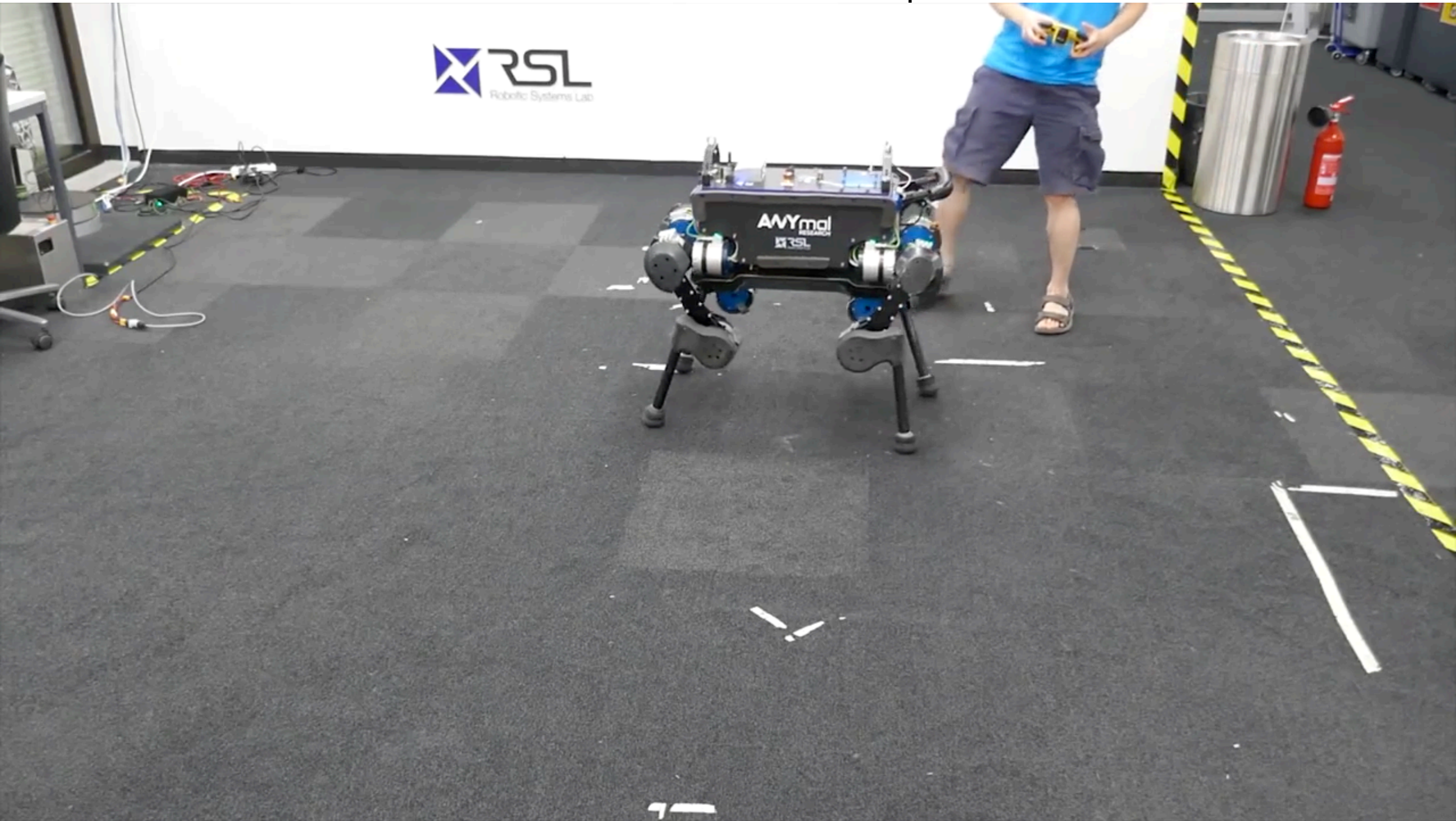
Hessel et. al Rainbow DQN, 2017



- Learning has been shown to be possible in simulation
- but can I use on a real robot??
- !!! millions (or billions) of real-world trials are needed



If exteroceptive sensors are not used, then transfer from accurate simulation is possible



[Hwangbo, ETH Zurich, Science Robotics, 2018]

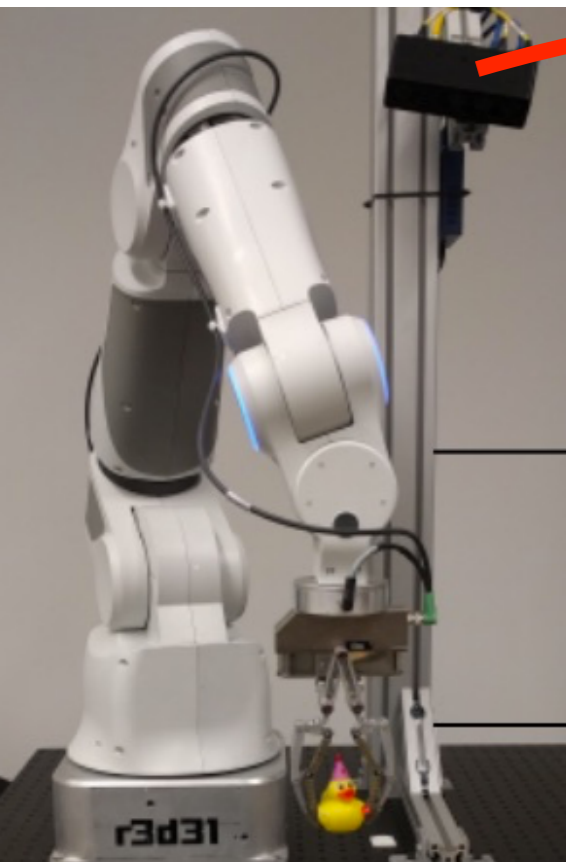


Czech Technical University in Prague
Faculty of Electrical Engineering, Department of Cybernetics

[Levine IJRR 2017] <https://arxiv.org/abs/1603.02199>

Another option is to avoid simulation completely !!!

manipulator+ RGB camera

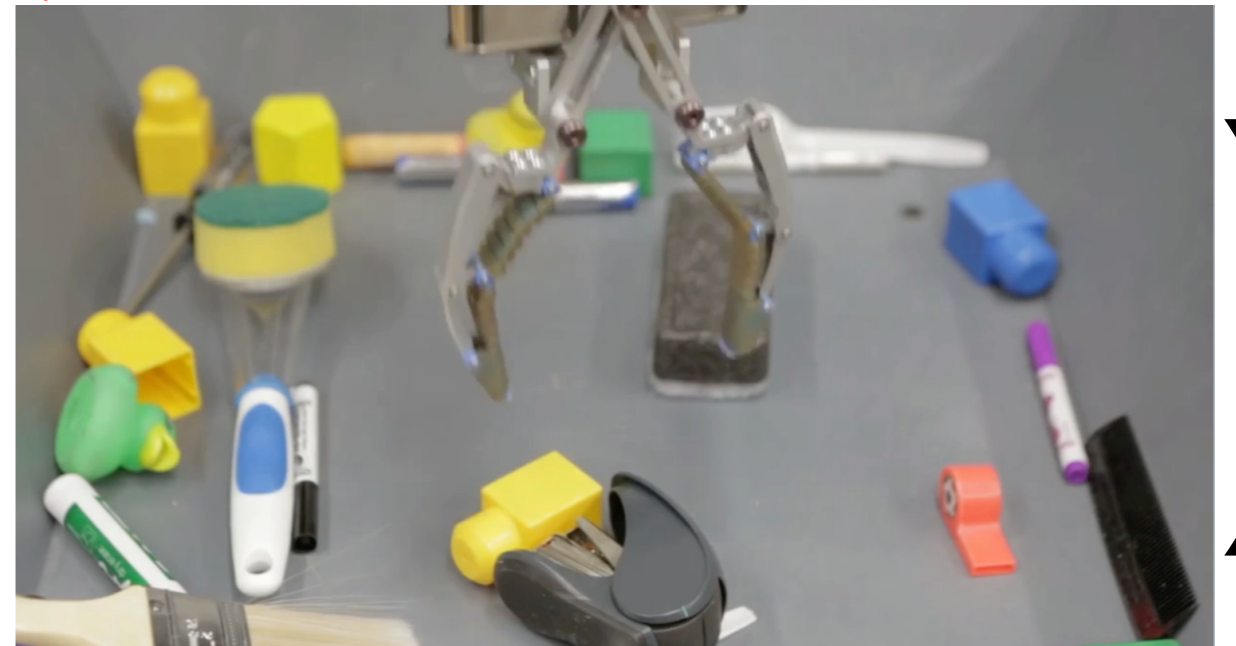


joint torques



image

$= \pi_{\theta}$



Continues motion control from RGB(D)



[Levine IJRR 2017] <https://arxiv.org/abs/1603.02199>

Source: Peter Pastor



Czech Technical University in Prague
Faculty of Electrical Engineering, Department of Cybernetics

Learning from expert demonstrations

- Sometimes easier to provide good trajectories than good rewards.



Learning from expert demonstrations

- Sometimes easier to provide good trajectories than good rewards.
- Imitation learning setup



Learning from expert demonstrations

- Sometimes easier to provide good trajectories than good rewards.
- Imitation learning setup
 1. Collect expert trajectories $\tau_1^*, \tau_2^*, \tau_3^*, \dots$
 2. Find policy $\arg \min_{\theta} \sum_{(\mathbf{x}_i, \mathbf{a}_i) \in \tau^*} \|\pi_{\theta}(\mathbf{x}_i) - \mathbf{a}_i\|_2^2$



Learning from expert demonstrations

- Sometimes easier to provide good trajectories than good rewards.
- Imitation learning setup (**statistically inconsistent+ blackbox**)
 1. Collect expert trajectories $\tau_1^*, \tau_2^*, \tau_3^*, \dots$
 2. Find policy $\arg \min_{\theta} \sum_{(\mathbf{x}_i, \mathbf{a}_i) \in \tau^*} \|\pi_{\theta}(\mathbf{x}_i) - \mathbf{a}_i\|_2^2$
- Inverse reinforcement learning setup



Learning from expert demonstrations

- Sometimes easier to provide good trajectories than good rewards.
- Imitation learning setup (statistically inconsistent+ blackbox)
 1. Collect expert trajectories $\tau_1^*, \tau_2^*, \tau_3^*, \dots$
 2. Find policy $\arg \min_{\theta} \sum_{(\mathbf{x}_i, \mathbf{a}_i) \in \tau^*} \|\pi_{\theta}(\mathbf{x}_i) - \mathbf{a}_i\|_2^2$
- Inverse reinforcement learning setup
 1. Collect expert trajectories $\tau_1^*, \tau_2^*, \tau_3^*, \dots$
 2. Find reward function $r_{\mathbf{w}}$

$$\arg \min_{\mathbf{w}} \|\mathbf{w}\|_2^2$$

$$\text{subject to: } \sum_{(\mathbf{x}, \mathbf{u}, \mathbf{x}') \in \tau^*} r_{\mathbf{w}}(\mathbf{x}, \mathbf{u}, \mathbf{x}') \leq \sum_{(\mathbf{x}, \mathbf{u}, \mathbf{x}') \in \{\mathcal{T} \setminus \tau^*\}} r_{\mathbf{w}}(\mathbf{x}, \mathbf{u}, \mathbf{x}')$$



Learning from expert demonstrations

- Sometimes easier to provide good trajectories than good rewards.
- Imitation learning setup (statistically inconsistent+ blackbox)
 1. Collect expert trajectories $\tau_1^*, \tau_2^*, \tau_3^*, \dots$
 2. Find policy $\arg \min_{\theta} \sum_{(\mathbf{x}_i, \mathbf{a}_i) \in \tau^*} \|\pi_{\theta}(\mathbf{x}_i) - \mathbf{a}_i\|_2^2$
- Inverse reinforcement learning setup
 1. Collect expert trajectories $\tau_1^*, \tau_2^*, \tau_3^*, \dots$
 2. Find reward function $r_{\mathbf{w}}$

$$\arg \min_{\mathbf{w}} \|\mathbf{w}\|_2^2$$

$$\text{subject to: } \sum_{(\mathbf{x}, \mathbf{u}, \mathbf{x}') \in \tau^*} r_{\mathbf{w}}(\mathbf{x}, \mathbf{u}, \mathbf{x}') \leq \sum_{(\mathbf{x}, \mathbf{u}, \mathbf{x}') \in \{\mathcal{T} \setminus \tau^*\}} r_{\mathbf{w}}(\mathbf{x}, \mathbf{u}, \mathbf{x}')$$

3. Solve underlying RL task



Abbeel et al. IJRR 2010

- inverse reinforcement learning
- **state space:** angular and euclidean position, velocity, acceleration
- **action space:** motor torques
- learning reward function from expert pilot



Abbeel et al. IJRR 2010



Silver et al. IJRR 2010



<http://www.dtic.mil/dtic/tr/fulltext/u2/a525288.pdf>

Czech Technical University in Prague

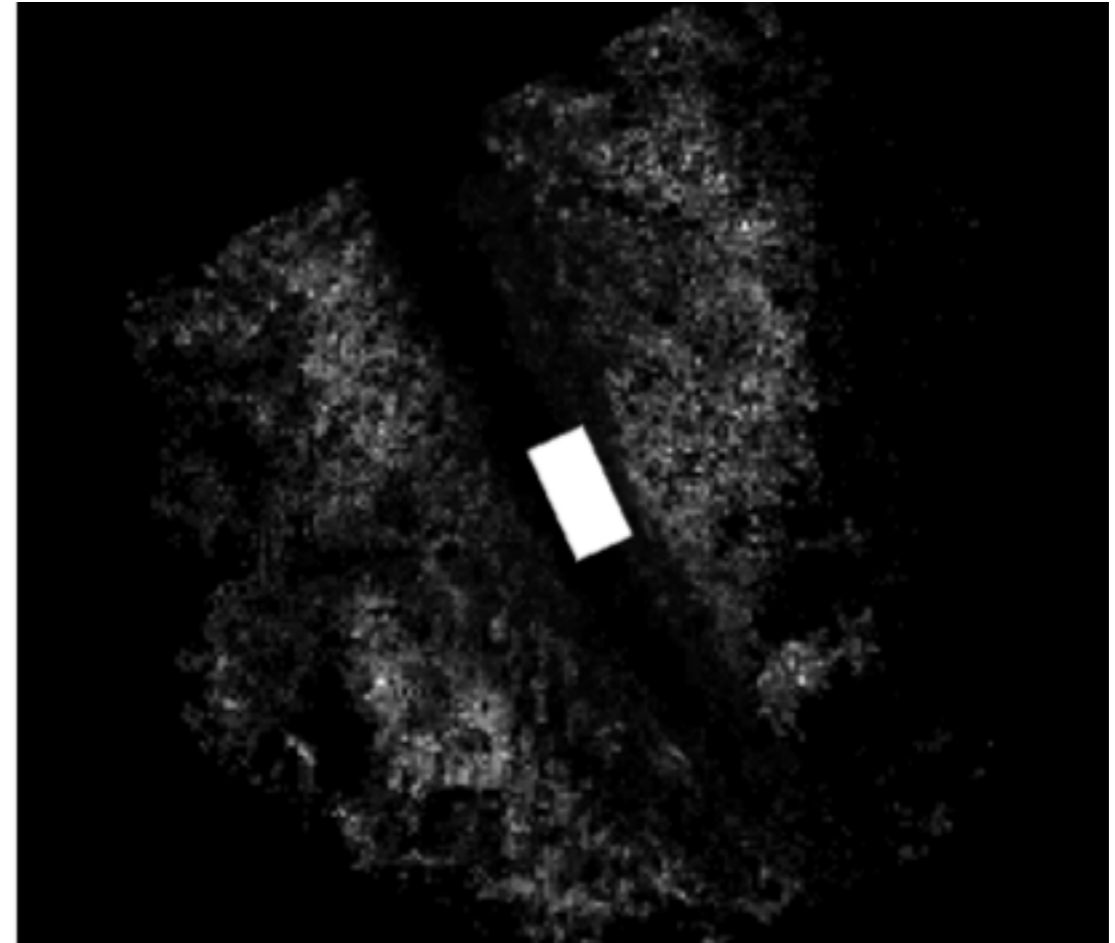
Faculty of Electrical Engineering, Department of Cybernetics



Silver et al. IJRR 2010



input image (state)



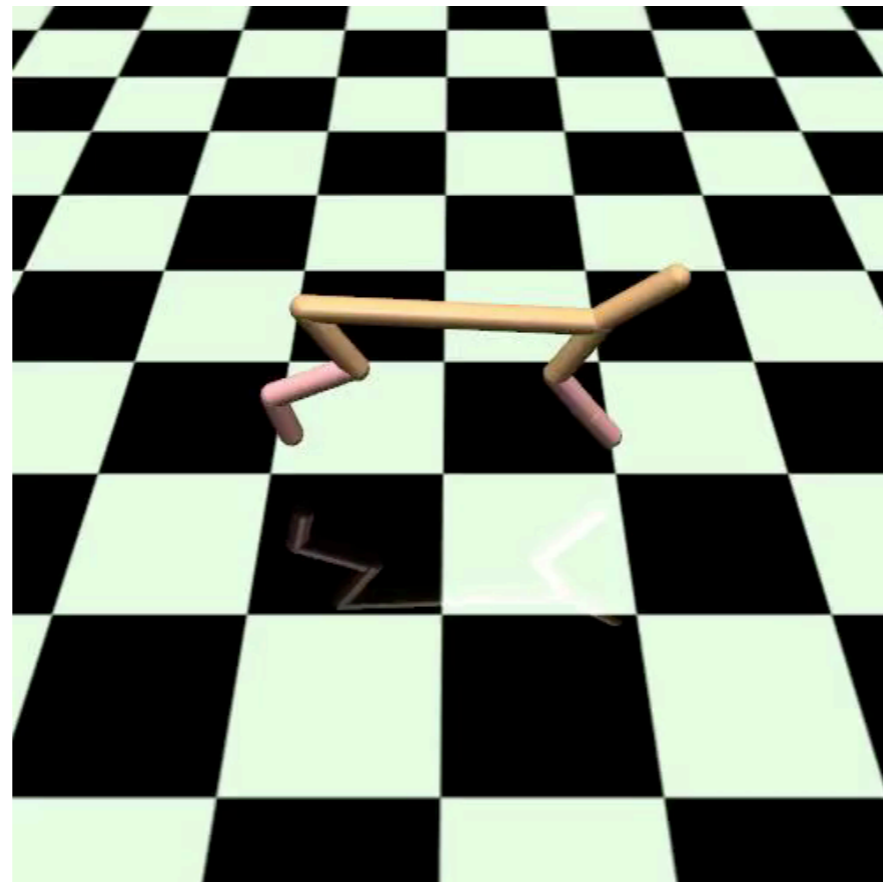
learned reward function
(traversability map)

<http://www.dtic.mil/dtic/tr/fulltext/u2/a525288.pdf>



Reward shaping

- Sparse rewards are easier to design correctly
- Dense rewards are easier to learn
- Half cheetah:
 - sparse rewards (for reaching the goal position fast)
 - dense rewards (for velocity)



Reward shaping

- Sparse rewards are easier to design correctly
- Dense rewards are easier to learn



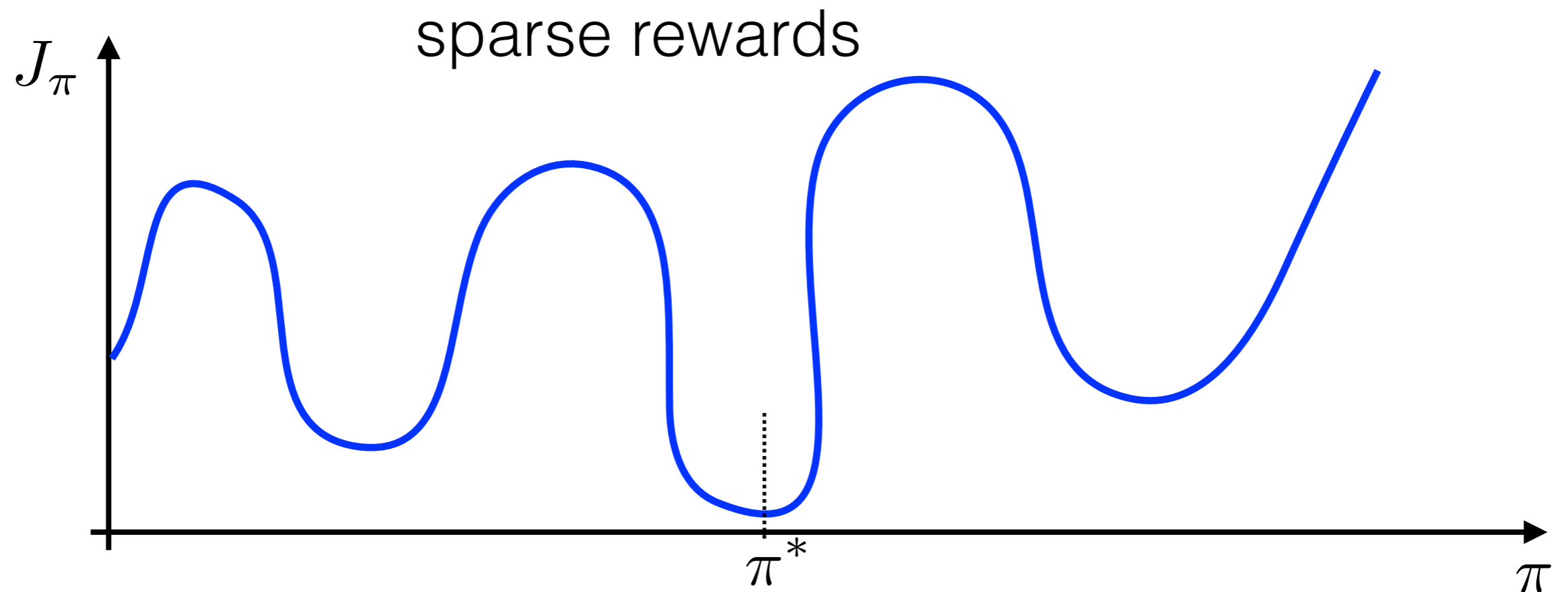
Reward shaping

- Sparse rewards are easier to design correctly
- Dense rewards are easier to learn



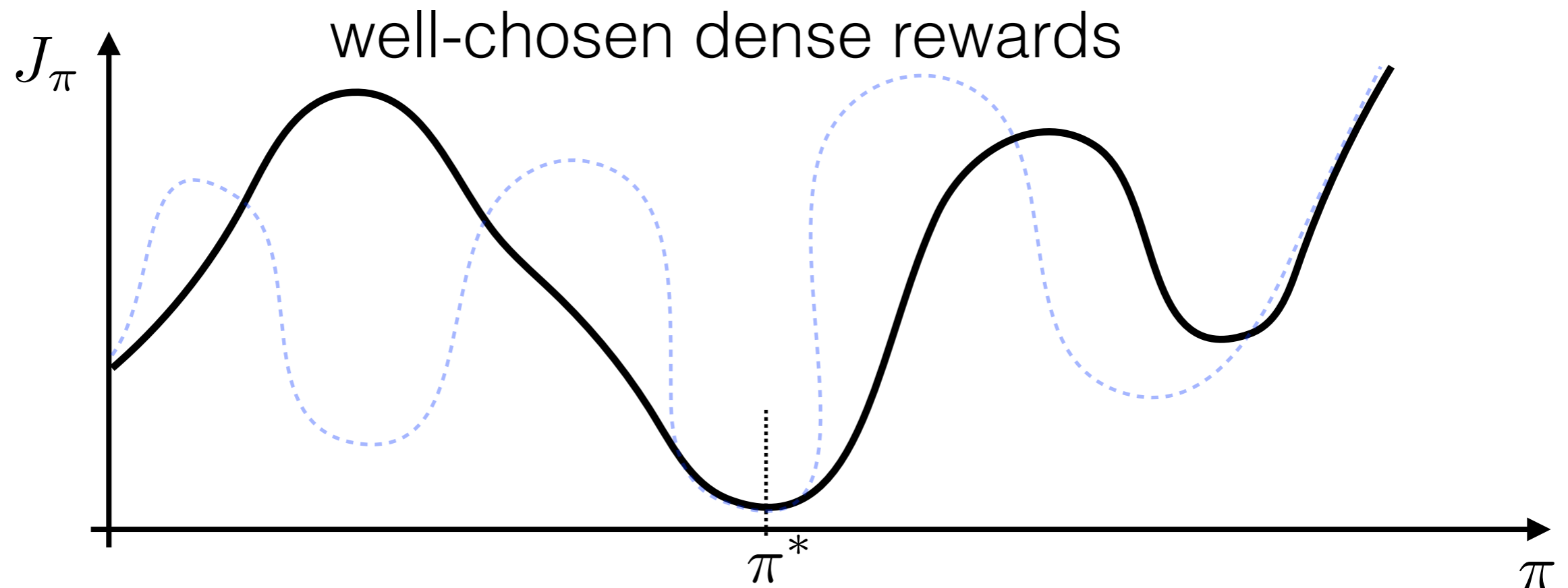
Reward shaping

- Sparse rewards are easier to design correctly
- Dense rewards are easier to learn



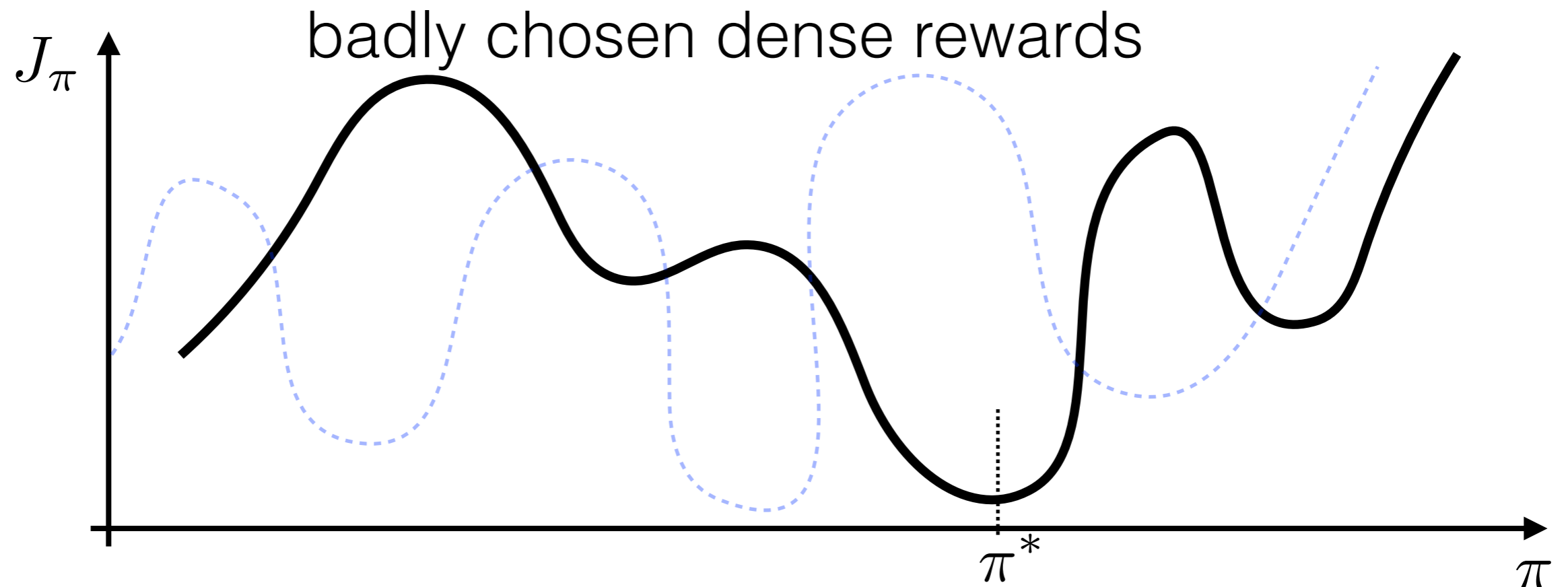
Reward shaping

- Sparse rewards are easier to design correctly
- Dense rewards are easier to learn



Reward shaping

- Sparse rewards are easier to design correctly
- Dense rewards are easier to learn



Reward shaping

- Sparse rewards are easier to design correctly
- Dense rewards are easier to learn
- Boat racing (bad dense rewards):
 - sparse rewards (winning the race)
 - dense rewards (collecting powerups, checkpoints ...)



Czech Technical University in Prague

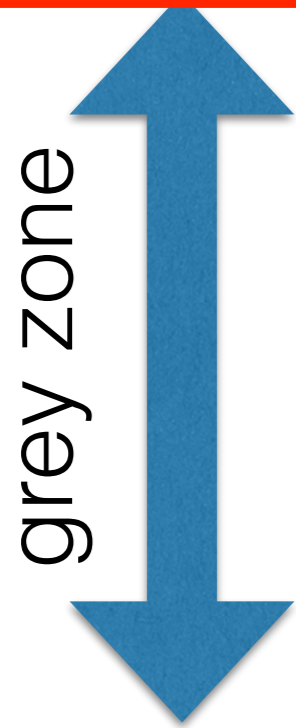
Faculty of Electrical Engineering, Department of Cybernetics



Taxonomy of policy search methods

- Direct policy search (primal task)

e.g. gradient ascent for $\pi^* = \arg \max_{\pi} J_{\pi}$



Episodic REPS [Peters, 2010]

PILCO [Deisenroth, ICML 2011]

Actor-critic (e.g. DPG [Silver, JMLR 2014])

Deep Q-learning (e.g. [Mnih, Nature 2015])

- Value-based methods (dual function [Kober, 2013])

e.g. search for $Q(\mathbf{x}, \mathbf{a}) = r(\mathbf{x}, \mathbf{a}, \mathbf{x}') + \gamma \max_{\mathbf{a}'} Q(\mathbf{x}', \mathbf{a}')$

$\pi^* = \arg \max_a Q(\mathbf{x}, \mathbf{a})$



Primal task

1. Randomly initialize policy π_θ



Primal task

1. Randomly initialize policy π_θ
2. Collect trajectories τ with policy π_θ



Primal task

1. Randomly initialize policy π_θ
2. Collect trajectories τ with policy π_θ
3. Denote $p(\tau|\pi_\theta)$ probability of τ occurs when following π_θ



Primal task

1. Randomly initialize policy π_θ
2. Collect trajectories τ with policy π_θ
3. Denote $p(\tau|\pi_\theta)$ probability of τ occurs when following π_θ
4. Define criterion

$$J(\theta) = \mathbb{E}_{\tau \sim p(\tau|\pi_\theta)} \{r(\tau)\} = \int_{\tau \in \mathcal{T}} p(\tau|\pi_\theta) r(\tau) d\tau \approx \frac{1}{N} \sum_{i=1}^N r(\tau_i)$$



Primal task

1. Randomly initialize policy π_θ
2. Collect trajectories τ with policy π_θ
3. Denote $p(\tau|\pi_\theta)$ probability of τ occurs when following π_θ
4. Define criterion

$$J(\theta) = \mathbb{E}_{\tau \sim p(\tau|\pi_\theta)} \{r(\tau)\} = \int_{\tau \in \mathcal{T}} p(\tau|\pi_\theta) r(\tau) d\tau \approx \frac{1}{N} \sum_{i=1}^N r(\tau_i)$$

5. Optimize criterion (e.g. gradient descent)

$$\theta^* = \arg \min_{\theta} J(\theta)$$

6. Repeat from 2



Primal task

$$J(\theta) = \mathbb{E}_{\tau \sim p(\tau|\pi_\theta)} \{r(\tau)\} = \int_{\tau \in \mathcal{T}} p(\tau|\pi_\theta) r(\tau) d\tau \approx \frac{1}{N} \sum_{i=1}^N r(\tau_i)$$

$$\theta^* = \arg \min_{\theta} J(\theta)$$

- What do I need for gradient descent optimization? $\frac{\partial J(\theta)}{\partial \theta}^\top$
- Perturb parameters by $\Delta\theta_i$ and estimate $J(\theta + \Delta\theta_i)$

$$J(\theta + \Delta\theta_i) = J(\theta) + \frac{\partial J(\theta)}{\partial \theta}^\top \Delta\theta_i$$
$$\Delta\theta_i^\top \frac{\partial J(\theta)}{\partial \theta} = J(\theta) - J(\theta + \Delta\theta_i)$$



Primal task

$$J(\theta) = \mathbb{E}_{\tau \sim p(\tau|\pi_\theta)} \{r(\tau)\} = \int_{\tau \in \mathcal{T}} p(\tau|\pi_\theta) r(\tau) d\tau \approx \frac{1}{N} \sum_{i=1}^N r(\tau_i)$$

$$\theta^* = \arg \min_{\theta} J(\theta)$$

- What do I need for gradient descent optimization? $\frac{\partial J(\theta)}{\partial \theta}^\top$
- Perturb parameters by $\Delta\theta_i$ and estimate $J(\theta + \Delta\theta_i)$

$$J(\theta + \Delta\theta_i) = J(\theta) + \frac{\partial J(\theta)}{\partial \theta}^\top \Delta\theta_i$$

$$\Delta\theta_i^\top \frac{\partial J(\theta)}{\partial \theta} = J(\theta) - J(\theta + \Delta\theta_i)$$

$$\underbrace{\begin{bmatrix} \Delta\theta_1^\top \\ \vdots \\ \Delta\theta_n^\top \end{bmatrix}}_{\text{matrix A}} \frac{\partial J(\theta)}{\partial \theta} = \underbrace{\begin{bmatrix} J(\theta) - J(\theta + \Delta\theta_1) \\ \vdots \\ J(\theta) - J(\theta + \Delta\theta_n) \end{bmatrix}}_{\text{vector b}}$$



Primal task

$$\underbrace{\begin{bmatrix} \Delta\theta_1^\top \\ \vdots \\ \Delta\theta_n^\top \end{bmatrix}}_{\text{matrix } \mathbf{A}} \frac{\partial J(\theta)}{\partial \theta} = \underbrace{\begin{bmatrix} J(\theta) - J(\theta + \Delta\theta_1) \\ \vdots \\ J(\theta) - J(\theta + \Delta\theta_n) \end{bmatrix}}_{\text{vector } \mathbf{b}}$$

$$\frac{\partial J(\theta)}{\partial \theta} = \begin{bmatrix} \Delta\theta_1^\top \\ \vdots \\ \Delta\theta_n^\top \end{bmatrix}^+ \cdot \begin{bmatrix} J(\theta) - J(\theta + \Delta\theta_1) \\ \vdots \\ J(\theta) - J(\theta + \Delta\theta_n) \end{bmatrix}$$



Primal task

1. Randomly initialize θ
2. Collect trajectories randomly perturbed policy $\pi_{\theta + \Delta\theta_i}$
3. Compute gradient $\frac{\partial J(\theta)}{\partial \theta}^\top$ using pseudo-inverse

$$\frac{\partial J(\theta)}{\partial \theta} = \begin{bmatrix} \Delta\theta_1^\top \\ \vdots \\ \Delta\theta_n^\top \end{bmatrix}^+ \cdot \begin{bmatrix} J(\theta) - J(\theta + \Delta\theta_1) \\ \vdots \\ J(\theta) - J(\theta + \Delta\theta_n) \end{bmatrix}$$

4. Update parameters

$$\theta \leftarrow \theta + \alpha \frac{\partial J(\theta)}{\partial \theta}$$



Primal task

REINFORCE: better gradient approximation

- stochastic policy

$$\pi_{\theta}(\mathbf{u}|\mathbf{x}) : X \times U \rightarrow [0; 1]$$

- gradient of the criterion

$$\nabla_{\theta} J(\theta) = \int_T \nabla_{\theta} p(\tau|\theta) r(\tau) d\tau$$

- likelihood ratio trick express gradient of the prob distr.

$$\nabla_{\theta} p(\tau|\theta) = p(\tau|\theta) \nabla_{\theta} \log p(\tau|\theta)$$



Primal task

- after substitution

$$\nabla_{\theta} J(\theta) = \int_{\mathcal{T}} p(\tau|\theta) \nabla_{\theta} \log p(\tau|\theta) r(\tau) d\tau =$$

$$= E[\nabla_{\theta} \log p(\tau|\theta) r(\tau)] \approx \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log p(\tau_i|\theta) r(\tau_i)$$

- where prob distribution simplified using MDP assumption

$$p(\tau|\theta) = p(\mathbf{x}_0) \prod_k p(\mathbf{x}_{k+1}|\mathbf{x}_k, \mathbf{u}_k) \pi_{\theta}(\mathbf{u}_k|\mathbf{x}_k)$$

$$\begin{aligned} \nabla_{\theta} \log p(\tau|\theta) &= \nabla_{\theta} [\log p(\mathbf{x}_0) + \sum_k \log p(\mathbf{x}_{k+1}|\mathbf{x}_k, \mathbf{u}_k) + \sum_k \log \pi_{\theta}(\mathbf{u}_k|\mathbf{x}_k)] \\ &= \sum_k \nabla_{\theta} \log \pi_{\theta}(\mathbf{u}_k|\mathbf{x}_k) \end{aligned}$$

Primal task

REINFORCE algorithm:

- collect N trajectories

$$\tau_1 = [(\mathbf{u}_{1,1}, \mathbf{x}_{1,1}) \dots \mathbf{u}_{M,1}, \mathbf{x}_{M,1}]$$

⋮

$$\tau_N = [(\mathbf{u}_{1,N}, \mathbf{x}_{1,N}) \dots \mathbf{u}_{M,N}, \mathbf{x}_{M,N}]$$

- compute gradient

$$\nabla_{\theta} J(\theta) = \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^M \nabla_{\theta} \log \pi_{\theta}(\mathbf{u}_{k,i} | \mathbf{x}_{k,i})$$

- update parameters

$$\theta \leftarrow \theta + \alpha \frac{\partial J(\theta)}{\partial \theta}$$



Primal task

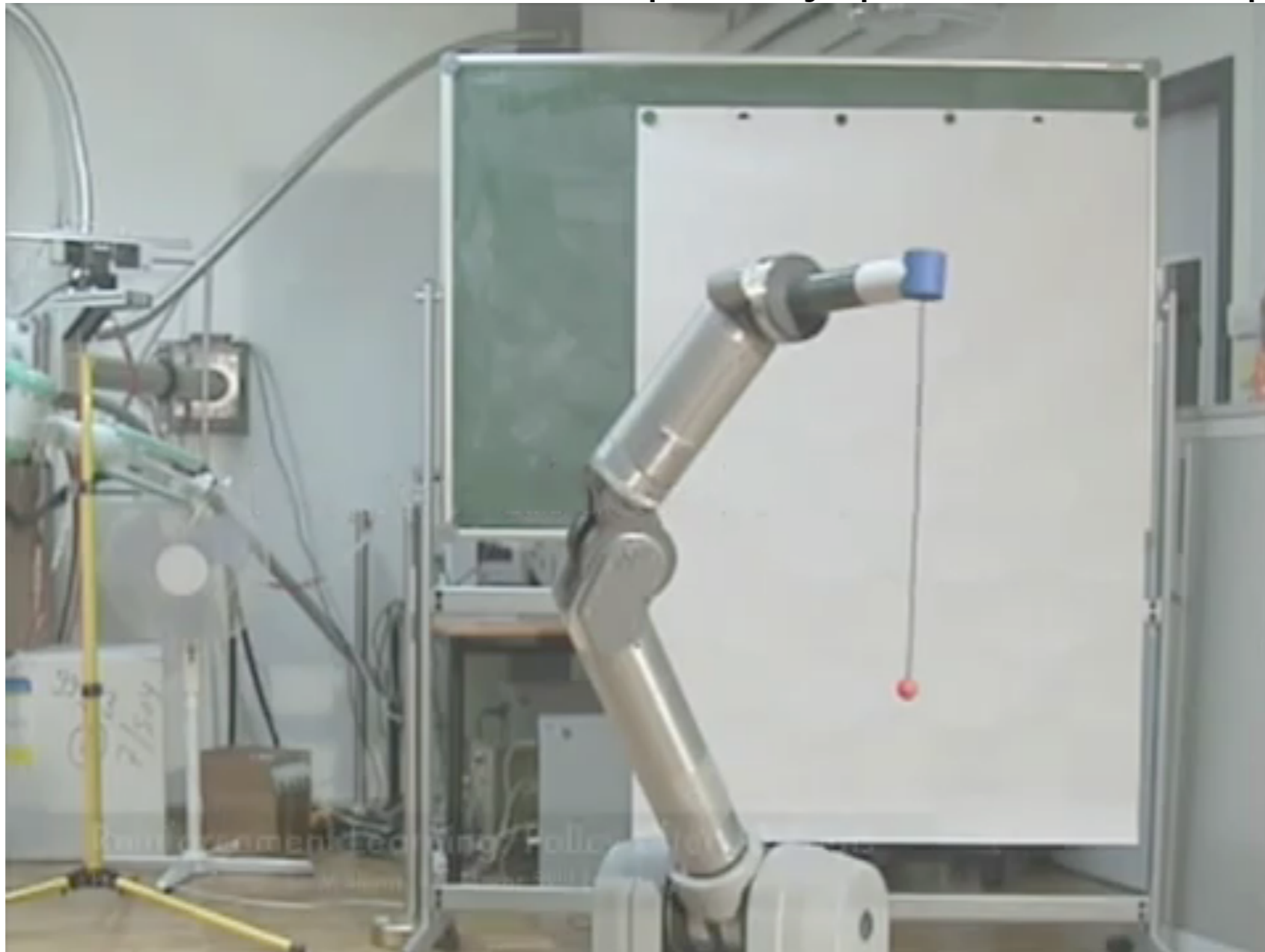
- No motion model required
- Converges to local optima (good initialization needed)
- High-dimensional parameters are requires many samples
- Imitation learning from expert trajectories

- There are better gradient approximations [Deisenroth 2013] (e.g. REINFORCE, GPREPS, ...)
[Deisenroth 2013] M. Deisenroth, G. Neumann and J. Peters, A Survey on Policy Search for Robotics, NOW, 2013



Peters et al. NOW 2013

- imitation learning from human demonstration
- **state space:** joint positions, velocities, acceler.
- **action space:** motor torques
- gradient minimization in policy parameter space



Summary RL

- No motion model required
- Converges to local optima (good initialization needed)
- High-dimensional parameters => requires many samples
- Imitation or Inverse RL learning from expert trajectories

- There are better gradient approximations [Deisenroth 2013] (e.g. REINFORCE, GPREPS, ...)
[Deisenroth 2013] M. Deisenroth, G. Neumann and J. Peters, A Survey on Policy Search for Robotics, NOW, 2013

- If motion model is available then trajectory optimization
[Tassa 2013] Tassa, Synthesis and Stabilization of Complex Behaviors through Online Trajectory Optimization, IROS2013



Taxonomy of policy search methods

- Direct policy search (primal task)

e.g. gradient ascent for $\pi^* = \arg \max_{\pi} J_{\pi}$

grey zone

Episodic REPS [Peters, 2010]

PILCO [Deisenroth, ICML 2011]

Actor-critic (e.g. DPG [Silver, JMLR 2014])

Deep Q-learning (e.g. [Mnih, Nature 2015])

- Value-based methods (dual function [Kober, 2013])

e.g. search for $Q(\mathbf{x}, \mathbf{a}) = r(\mathbf{x}, \mathbf{a}, \mathbf{x}') + \gamma \max_{\mathbf{a}'} Q(\mathbf{x}', \mathbf{a}')$

$\pi^* = \arg \max_a Q(\mathbf{x}, \mathbf{a})$



Actor-critic methods

1. Collect trajectories $\tau_1, \tau_2, \tau_3, \dots$. initialize $\theta = \text{rand}$
2. Estimate $y = r(\mathbf{x}, \mathbf{u}, \mathbf{x}') + \gamma \max_{\mathbf{u}'} Q_{\theta}(\mathbf{x}', \mathbf{u}')$
3. Update parameters by learning \mathbf{u}'

$$\arg \min_{\theta} \sum_{\mathbf{x}, \mathbf{u}, \mathbf{y}} \|Q_{\theta}(\mathbf{x}, \mathbf{u}) - \mathbf{y}\|$$

Approximated Q-learning



Actor-critic methods

1. Collect trajectories $\tau_1, \tau_2, \tau_3, \dots$. initialize $\theta = \text{rand}$
2. Estimate $y = r(\mathbf{x}, \mathbf{u}, \mathbf{x}') + \gamma \max_{\mathbf{u}'} Q_{\theta}(\mathbf{x}', \mathbf{u}')$
3. Update parameters by learning \mathbf{u}'

$$\arg \min_{\theta} \sum_{\mathbf{x}, \mathbf{u}, \mathbf{y}} \|Q_{\theta}(\mathbf{x}, \mathbf{u}) - \mathbf{y}\|$$

4. Learn policy π_{ω} which do actions maximizing the state-action value function on the collected trajectories

$$\arg \max_{\omega} \sum_{\mathbf{x} \in \tau} Q_{\theta}(\mathbf{x}, \pi_{\omega}(\mathbf{x}))$$

Direct policy optimization on Q



3D humanoid

Degrees-of-freedom: **22**

- 6 spatial
- 2 abdomen
- 2·2 shoulders
- 2·1 elbows
- 2·2 hips
- 2·1 knees
- 2·1 ankles

Control dimensions: **16** all joints

Cost:

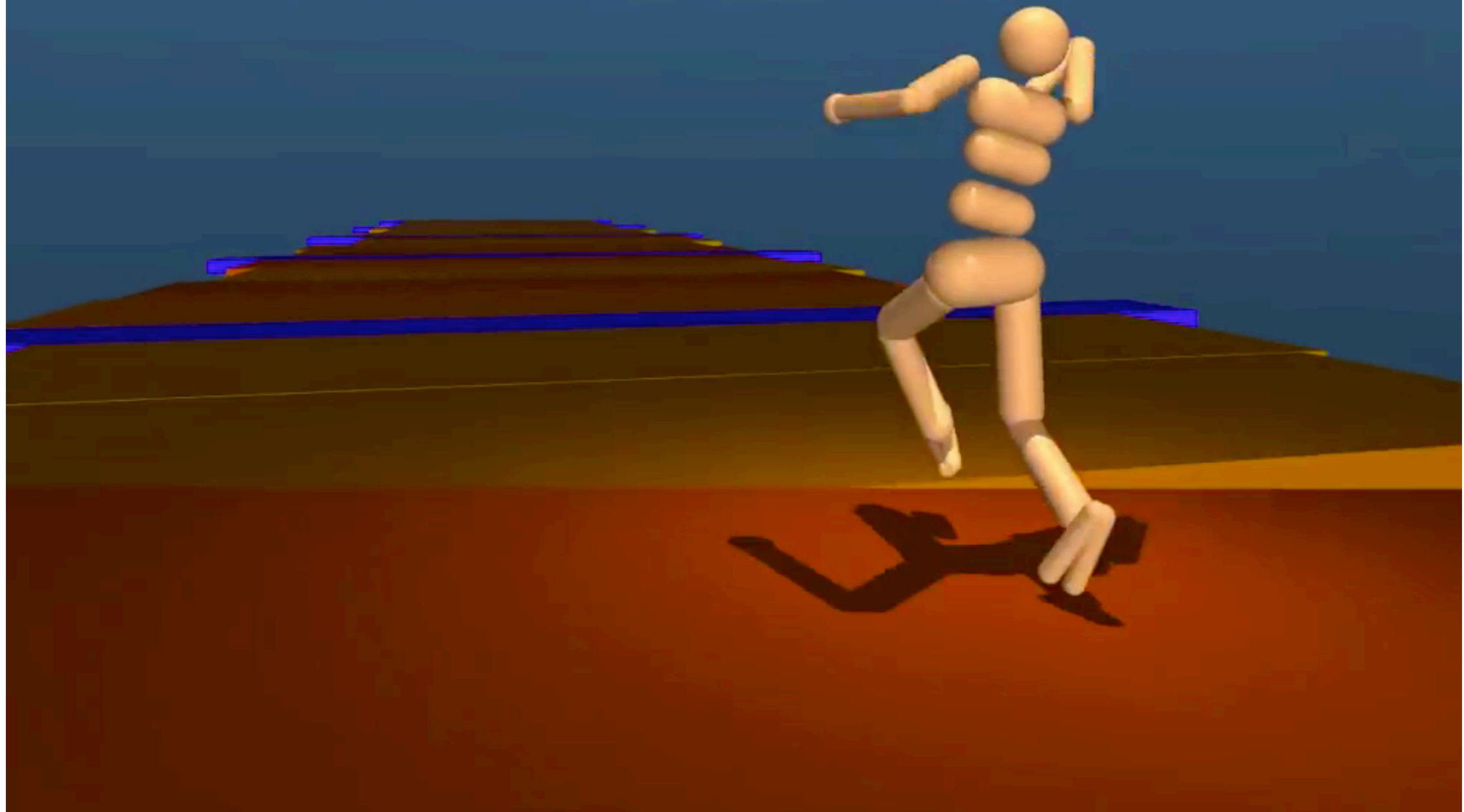
CoM over mean of feet, (in xy) + torso over CoM (in xy) + torso 1.3m over mean of feet (in z)

+ minimize horizontal torso velocity + minimize actuation



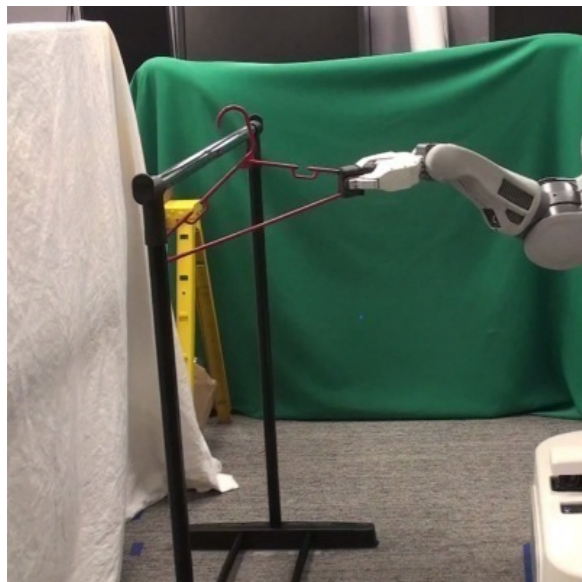
[Heess 2017] <https://arxiv.org/abs/1707.02286>

This agent, trained on several terrain types, has never seen the "see-saw" terrain.



Levine et al JMLR 2016

- guides policy gradient method by optimal trajectories
- **state space:** RGB camera images
- **action space:** motor torques



(a) hanger



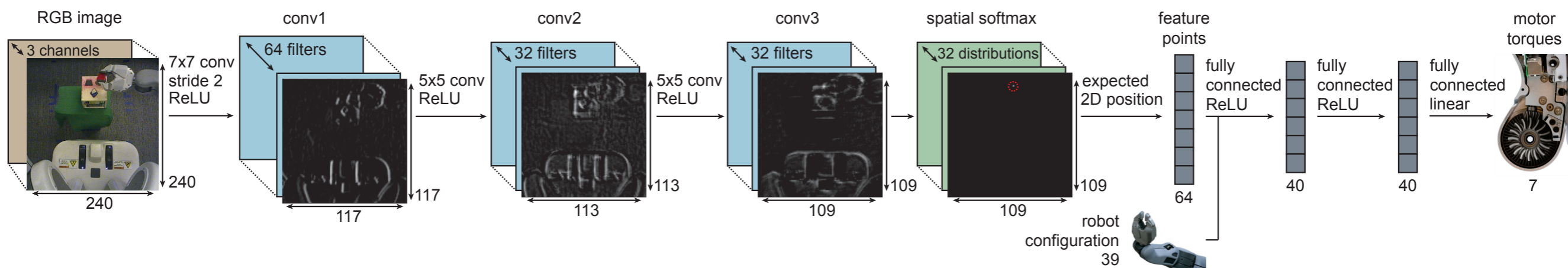
(b) cube



(c) hammer



(d) bottle



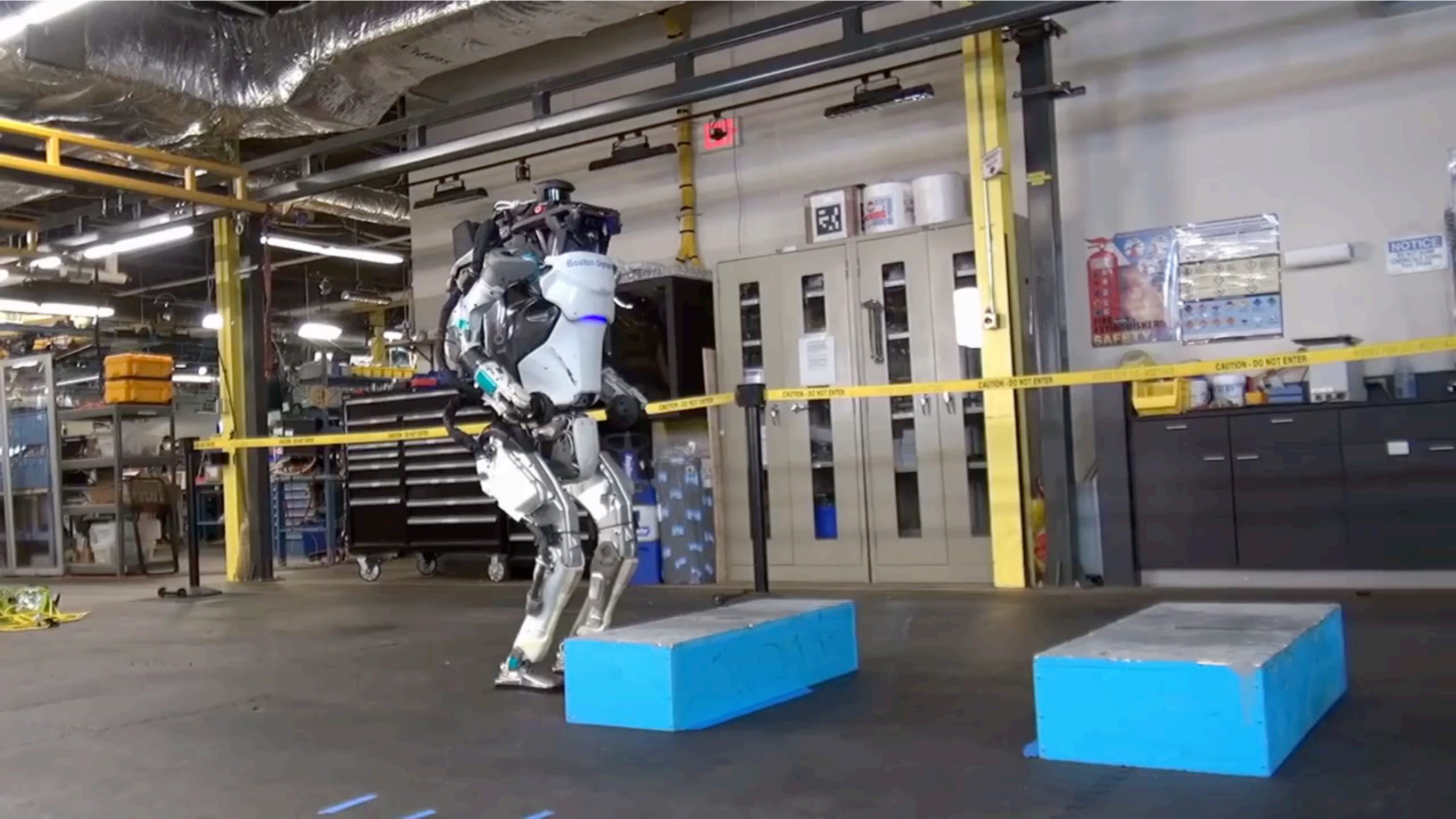
Levine et al JMLR 2016

Learned Visuomotor Policy: Bottle Task



Czech Technical University in Prague
Faculty of Electrical Engineering, Department of Cybernetics

Boston dynamics - Atlas - NO RL AT ALL



Boston dynamics - Big dog - NO RL AT ALL



Known RL successes

- AlphaGo/Alpha Zero <https://en.wikipedia.org/wiki/AlphaZero>
- SearchTrees has no chance in huge state-action spaces
 - AlphaGo:
 - beat professional Go player
 - 9 dan professional ranking
 - Alpha Zero: Top Chess Engine Championship 2017
 - 9h of self-play, no openingbooks nor endgames tables
 - 1 minute per move, 1GB RAM
 - 28 wins, 72 withdraws
- DOTA 2 openAI+ bot <https://blog.openai.com/dota-2/>
- AutoML <https://cloud.google.com/automl/>
 - [Zoph 2016] REINFORCE learns RCNN policy which generates deep CNN architectures.



Summary

- If accurate differentiable motion model and reward functions are known, than optimal control in MDP is straightforward optimization problem (efficiently tackled by DP or DDP)
- State-action value function is dual variable wrt policy. It serves as auxiliary function in the policy optimization:
 - actor-critic methods
 - heuristic in planning methods (LQR trees)
- **Holy grail** is to efficiently combine motion model, state-action value function with efficient planning, learning and exploration.
- RL will be much more useful for motion control, when accurate domain transfer methods (from simulators to reality) become available.

